

Part II: Sparse Gaussian Processes

Daniel Hernández-Lobato

Computer Science Department
Universidad Autónoma de Madrid

<http://dhnz1.org>, daniel.hernandez@uam.es

Computational Cost of Gaussian Processes

The memory cost is in $\mathcal{O}(N^2)$ since we have to compute Σ .

Computational Cost of Gaussian Processes

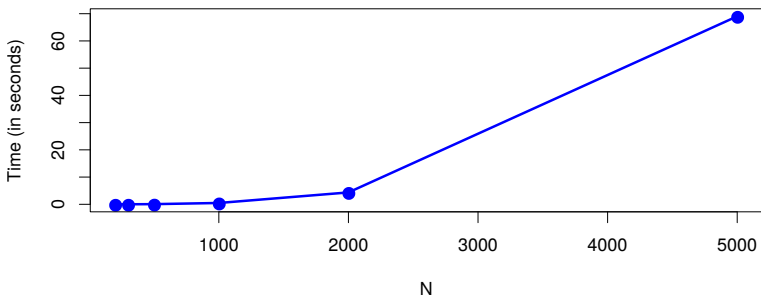
The memory cost is in $\mathcal{O}(N^2)$ since we have to compute Σ .

The computational cost is in $\mathcal{O}(N^3)$ since we have to invert Σ .

Computational Cost of Gaussian Processes

The memory cost is in $\mathcal{O}(N^2)$ since we have to compute Σ .

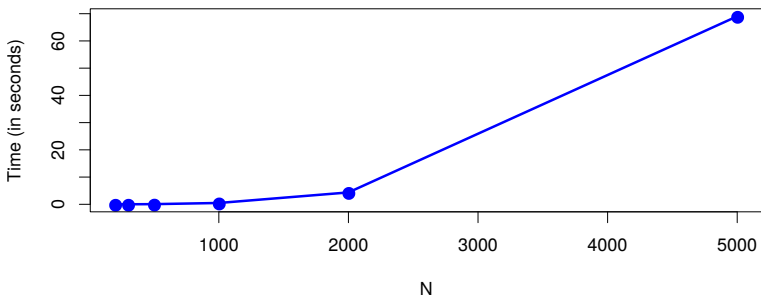
The computational cost is in $\mathcal{O}(N^3)$ since we have to invert Σ .



Computational Cost of Gaussian Processes

The memory cost is in $\mathcal{O}(N^2)$ since we have to compute Σ .

The computational cost is in $\mathcal{O}(N^3)$ since we have to invert Σ .



We can handle just a few thousand data instances at most!

Run the first cells of the notebook to measure GP fitting time and complete task 1!

Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

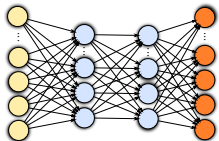
Idea: go back to the parametric model, but in such a way that we can still make inference easily!

Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

Idea: go back to the parametric model, but in such a way that we can still make inference easily!



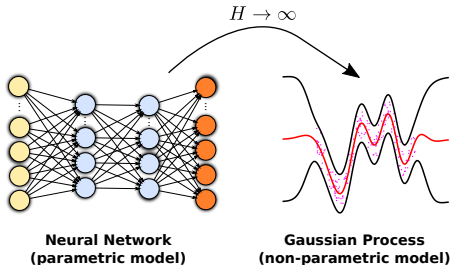
Neural Network
(parametric model)

Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

Idea: go back to the parametric model, but in such a way that we can still make inference easily!

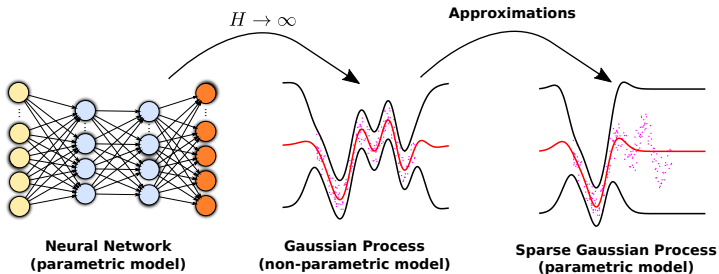


Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

Idea: go back to the parametric model, but in such a way that we can still make inference easily!

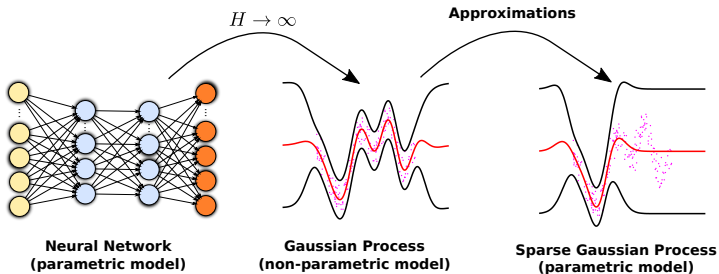


Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with N !

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

Idea: go back to the parametric model, but in such a way that we can still make inference easily!



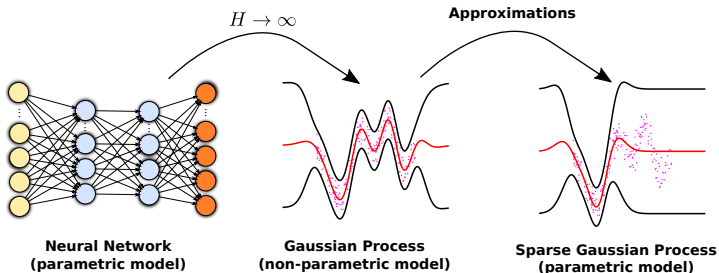
- **Nyström, Random Features and FITC:** approximate GP prior!

Improving the Cost of Gaussian Processes

GPs are non-parametric models whose flexibility grows with $N!$

GPs are the limiting case ($H \rightarrow \infty$) of Bayesian Neural Networks!

Idea: go back to the parametric model, but in such a way that we can still make inference easily!



- **Nyström, Random Features and FITC:** approximate GP prior!
- **VFE:** does approximate inference with a simplified distribution q .

The Nyström Method

Motivation: The posterior mean and covariances require $(\mathbf{I}\sigma^2 + \mathbf{\Sigma})^{-1}$.

The Nyström Method

Motivation: The posterior mean and covariances require $(\mathbf{I}\sigma^2 + \mathbf{\Sigma})^{-1}$.

Can we approximate the inverse of $\mathbf{I}\sigma^2 + \mathbf{\Sigma}$ with a cheaper cost?

The Nyström Method

Motivation: The posterior mean and covariances require $(\mathbf{I}\sigma^2 + \Sigma)^{-1}$.

Can we approximate the inverse of $\mathbf{I}\sigma^2 + \Sigma$ with a cheaper cost?

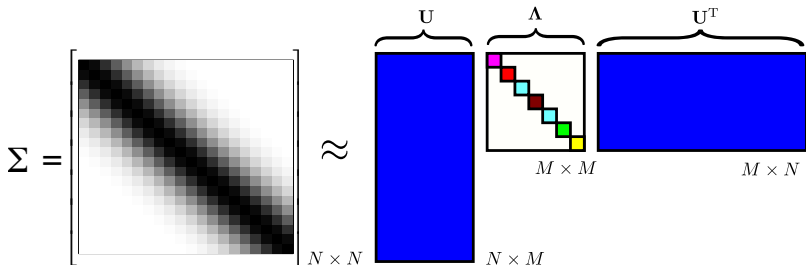
A low rank m approximation of Σ does the job:

The Nyström Method

Motivation: The posterior mean and covariances require $(\mathbf{I}\sigma^2 + \Sigma)^{-1}$.

Can we approximate the inverse of $\mathbf{I}\sigma^2 + \Sigma$ with a cheaper cost?

A low rank m approximation of Σ does the job:



The Nyström Method

Motivation: The posterior mean and covariances require $(\mathbf{I}\sigma^2 + \Sigma)^{-1}$.

Can we approximate the inverse of $\mathbf{I}\sigma^2 + \Sigma$ with a cheaper cost?

A low rank m approximation of Σ does the job:

$$\Sigma = \begin{bmatrix} \text{[Grayscale Heatmap]} \end{bmatrix}_{N \times N} \approx \underbrace{\begin{bmatrix} \text{[Blue Matrix]} \end{bmatrix}}_U_{N \times M} \underbrace{\begin{bmatrix} \text{[Diagonal Matrix]} \end{bmatrix}}_{\Lambda}_{M \times M} \underbrace{\begin{bmatrix} \text{[Blue Matrix]} \end{bmatrix}}_{U^T}_{M \times N}$$

The Woodbury formula gives $(\mathbf{I}\sigma^2 + \mathbf{U}\Lambda\mathbf{U}^T)^{-1}$ with cost $\mathcal{O}(M^2N)$!

Woodbury Formula

$$(\mathbf{A} + \mathbf{P}\mathbf{C}\mathbf{Q})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

Woodbury Formula

$$(\mathbf{A} + \mathbf{P}\mathbf{C}\mathbf{Q})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

Let us now use $\mathbf{A} = \mathbf{I}\sigma^2$, $\mathbf{P} = \mathbf{U}$, $\mathbf{Q} = \mathbf{U}^T$ and $\mathbf{C} = \mathbf{\Lambda}$.

Woodbury Formula

$$(\mathbf{A} + \mathbf{P}\mathbf{C}\mathbf{Q})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

Let us now use $\mathbf{A} = \mathbf{I}\sigma^2$, $\mathbf{P} = \mathbf{U}$, $\mathbf{Q} = \mathbf{U}^T$ and $\mathbf{C} = \mathbf{\Lambda}$.

Note that \mathbf{A} and \mathbf{C} are diagonal with sizes $N \times N$ and $M \times M$!

Woodbury Formula

$$(\mathbf{A} + \mathbf{P}\mathbf{C}\mathbf{Q})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

Let us now use $\mathbf{A} = \mathbf{I}\sigma^2$, $\mathbf{P} = \mathbf{U}$, $\mathbf{Q} = \mathbf{U}^T$ and $\mathbf{C} = \mathbf{\Lambda}$.

Note that \mathbf{A} and \mathbf{C} are diagonal with sizes $N \times N$ and $M \times M$!

$\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P} = \mathbf{\Lambda}^{-1} + \mathbf{U}^T\mathbf{U}\sigma^{-2}$ has size $M \times M$!

Woodbury Formula

$$(\mathbf{A} + \mathbf{P}\mathbf{C}\mathbf{Q})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

Let us now use $\mathbf{A} = \mathbf{I}\sigma^2$, $\mathbf{P} = \mathbf{U}$, $\mathbf{Q} = \mathbf{U}^T$ and $\mathbf{C} = \mathbf{\Lambda}$.

Note that \mathbf{A} and \mathbf{C} are diagonal with sizes $N \times N$ and $M \times M$!

$\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P} = \mathbf{\Lambda}^{-1} + \mathbf{U}^T\mathbf{U}\sigma^{-2}$ has size $M \times M$!

$$(\mathbf{I}\sigma^2 + \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T)^{-1} = \mathbf{I}\sigma^{-2} - \sigma^{-2}\mathbf{U}(\mathbf{\Lambda}^{-1} + \mathbf{U}^T\sigma^{-2}\mathbf{U})^{-1}\mathbf{U}^T\sigma^{-2}$$

Woodbury Formula

$$(\mathbf{A} + \mathbf{P}\mathbf{C}\mathbf{Q})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

Let us now use $\mathbf{A} = \mathbf{I}\sigma^2$, $\mathbf{P} = \mathbf{U}$, $\mathbf{Q} = \mathbf{U}^T$ and $\mathbf{C} = \mathbf{\Lambda}$.

Note that \mathbf{A} and \mathbf{C} are diagonal with sizes $N \times N$ and $M \times M$!

$\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P} = \mathbf{\Lambda}^{-1} + \mathbf{U}^T\mathbf{U}\sigma^{-2}$ has size $M \times M$!

$$(\mathbf{I}\sigma^2 + \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T)^{-1} = \mathbf{I}\sigma^{-2} - \sigma^{-2}\mathbf{U}(\mathbf{\Lambda}^{-1} + \mathbf{U}^T\sigma^{-2}\mathbf{U})^{-1}\mathbf{U}^T\sigma^{-2}$$

Computing the whole expression has cost $\mathcal{O}(NM^2)$!

Eigenfunction Analysis of Covariance Functions

GPs are equivalent to a Bayesian linear model on an extended input space given by the eigenfunctions of the covariance function.

Eigenfunction Analysis of Covariance Functions

GPs are equivalent to a Bayesian linear model on an extended input space given by the eigenfunctions of the covariance function.

Extended input space: A function $\phi(\cdot)$ that obeys

$$\int C(\mathbf{x}, \mathbf{x}')\phi(\mathbf{x})p(\mathbf{x})d\mathbf{x} = \lambda\phi(\mathbf{x}'),$$

is an eigenfunction of $C(\cdot, \cdot)$ with eigenvalue λ , w.r.t., $p(\mathbf{x})$.

Eigenfunction Analysis of Covariance Functions

GPs are equivalent to a Bayesian linear model on an extended input space given by the eigenfunctions of the covariance function.

Extended input space: A function $\phi(\cdot)$ that obeys

$$\int C(\mathbf{x}, \mathbf{x}') \phi(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \lambda \phi(\mathbf{x}'),$$

is an eigenfunction of $C(\cdot, \cdot)$ with eigenvalue λ , w.r.t., $p(\mathbf{x})$.

Mercer's theorem:

$$C(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}').$$

An Analytic Example

Consider:

$$p(x) = \mathcal{N}(x|0, \sigma^2), \quad C(x, x') = \exp \left\{ -\frac{1}{2\ell^2} (x - x')^2 \right\} .$$

An Analytic Example

Consider:

$$p(x) = \mathcal{N}(x|0, \sigma^2), \quad C(x, x') = \exp \left\{ -\frac{1}{2\ell^2} (x - x')^2 \right\}.$$

Then,

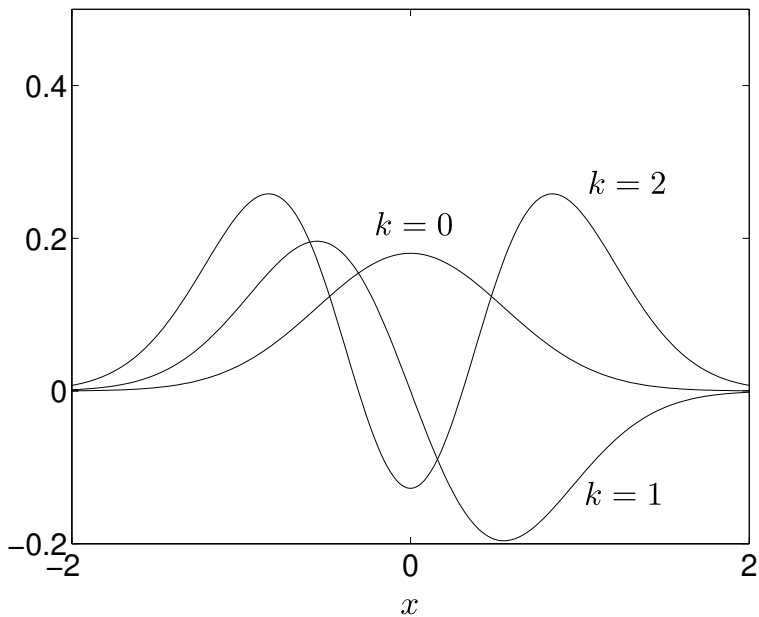
$$\lambda_k = \sqrt{\frac{2a}{A}} B^k, \quad \phi_k(x) = \exp \{ -(c - a)x^2 \} H_k(\sqrt{2c}x),$$

for $k = 0, 1, 2, \dots$, with

$$a^{-1} = 4\sigma^2, \quad b^{-1} = 2\ell^2, \quad c = \sqrt{a^2 + 2ab}, \quad A = a + b + c, \quad B = b/A,$$

and $H_k(\cdot)$, the k -th order Hermite polynomial.

Hermite Polynomials

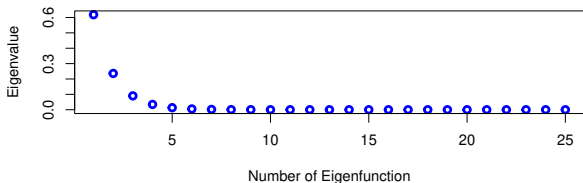


Covariance Function Approximation

Considering only the first eigenfunctions and eigenvalues is expected to give a good approximation of the covariance function!

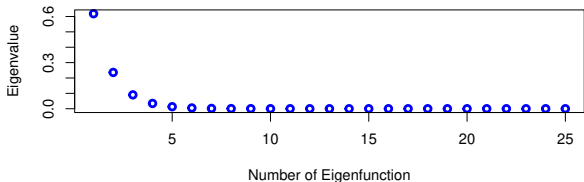
Covariance Function Approximation

Considering only the first eigenfunctions and eigenvalues is expected to give a good approximation of the covariance function!

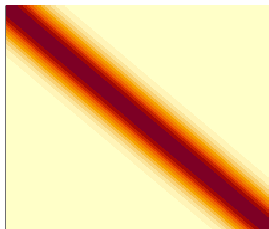


Covariance Function Approximation

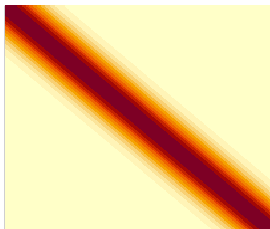
Considering only the first eigenfunctions and eigenvalues is expected to give a good approximation of the covariance function!



Exact Covariance Matrix



Approx. Covariance Matrix



Nyström Approximation of Eigenfunctions

Let $p(x)$ be the distribution of the observed data.

Nyström Approximation of Eigenfunctions

Let $p(\mathbf{x})$ be the distribution of the observed data.

Consider the Monte Carlo estimator:

$$\lambda_i \phi_i(\mathbf{x}') = \int C(\mathbf{x}, \mathbf{x}') \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N} \sum_{n=1}^N C(\mathbf{x}_n, \mathbf{x}') \phi_i(\mathbf{x}_n).$$

Nyström Approximation of Eigenfunctions

Let $p(\mathbf{x})$ be the distribution of the observed data.

Consider the Monte Carlo estimator:

$$\lambda_i \phi_i(\mathbf{x}') = \int C(\mathbf{x}, \mathbf{x}') \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N} \sum_{n=1}^N C(\mathbf{x}_n, \mathbf{x}') \phi_i(\mathbf{x}_n).$$

For \mathbf{x}' in the training set, this motivates the following eigenvalue problem:

$$\lambda_i^{\text{mat}} \mathbf{u}_i = \Sigma \mathbf{u}_i,$$

with $\Sigma_{i,j} = C(\mathbf{x}_i, \mathbf{x}_j)$. Then, we approximate $\phi_i(\mathbf{x}_j) \approx \sqrt{N}(\mathbf{u}_i)_j = \tilde{\phi}_i(\mathbf{x}_j)$, and $\lambda_i \approx \lambda_i^{\text{mat}}/N = \tilde{\lambda}_i$, which guarantees that $\Sigma = \tilde{\Phi} \tilde{\Lambda} \tilde{\Phi}^T$, with $j = 1, \dots, N$.

Nyström Approximation of Eigenfunctions

Let $p(\mathbf{x})$ be the distribution of the observed data.

Consider the Monte Carlo estimator:

$$\lambda_i \phi_i(\mathbf{x}') = \int C(\mathbf{x}, \mathbf{x}') \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N} \sum_{n=1}^N C(\mathbf{x}_n, \mathbf{x}') \phi_i(\mathbf{x}_n).$$

For \mathbf{x}' in the training set, this motivates the following eigenvalue problem:

$$\lambda_i^{\text{mat}} \mathbf{u}_i = \Sigma \mathbf{u}_i,$$

with $\Sigma_{i,j} = C(\mathbf{x}_i, \mathbf{x}_j)$. Then, we approximate $\phi_i(\mathbf{x}_j) \approx \sqrt{N}(\mathbf{u}_i)_j = \tilde{\phi}_i(\mathbf{x}_j)$, and $\lambda_i \approx \lambda_i^{\text{mat}}/N = \tilde{\lambda}_i$, which guarantees that $\Sigma = \tilde{\Phi} \tilde{\Lambda} \tilde{\Phi}^T$, with $j = 1, \dots, N$.

For an arbitrary \mathbf{x}' not in the training set, then:

$$\tilde{\phi}_i(\mathbf{x}') = \frac{1}{N\lambda_i} \sum_{n=1}^N C(\mathbf{x}', \mathbf{x}_n) \phi_i(\mathbf{x}_n) \approx \frac{\sqrt{N}}{\lambda_i^{\text{mat}}} \sum_{n=1}^N C(\mathbf{x}', \mathbf{x}_n) (\mathbf{u}_i)_n = \frac{\sqrt{N}}{\lambda_i^{\text{mat}}} \Sigma(\mathbf{x}')^T \mathbf{u}_i.$$

Putting All Together

We choose a random subset of size $M < N$ of the training data, to approximate the eigenfunctions and eigenvalues!

Putting All Together

We choose a random subset of size $M < N$ of the training data, to approximate the eigenfunctions and eigenvalues!

Using Mercer's theorem and the previous approximation, we approximate the covariance function as:

$$C(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}') \approx \sum_{i=1}^M \tilde{\lambda}_i \tilde{\phi}_i(\mathbf{x}) \tilde{\phi}_i(\mathbf{x}').$$

Putting All Together

We choose a random subset of size $M < N$ of the training data, to approximate the eigenfunctions and eigenvalues!

Using Mercer's theorem and the previous approximation, we approximate the covariance function as:

$$C(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}') \approx \sum_{i=1}^M \tilde{\lambda}_i \tilde{\phi}_i(\mathbf{x}) \tilde{\phi}_i(\mathbf{x}').$$

which results in a rank M approximation of the covariance matrix Σ :

$$\Sigma \approx \tilde{\Sigma} = \tilde{\Phi} \tilde{\Lambda} \tilde{\Phi}^T = \Sigma_{N,M} \Sigma_{M,M}^{-1} \Sigma_{M,N}.$$

Putting All Together

We choose a random subset of size $M < N$ of the training data, to approximate the eigenfunctions and eigenvalues!

Using Mercer's theorem and the previous approximation, we approximate the covariance function as:

$$C(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{x}') \approx \sum_{i=1}^M \tilde{\lambda}_i \tilde{\phi}_i(\mathbf{x}) \tilde{\phi}_i(\mathbf{x}').$$

which results in a rank M approximation of the covariance matrix Σ :

$$\Sigma \approx \tilde{\Sigma} = \tilde{\Phi} \tilde{\Lambda} \tilde{\Phi}^T = \Sigma_{N,M} \Sigma_{M,M}^{-1} \Sigma_{M,N}.$$

The inverse of $I\sigma^2 + \tilde{\Sigma}$ can be efficiently computed using the Woodbury formula with cost $\mathcal{O}(NM^2)$!

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

Predictive Distribution

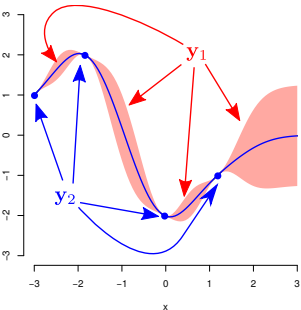
We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \Sigma_{\mathbf{f}\mathbf{f}^*} \\ \Sigma_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

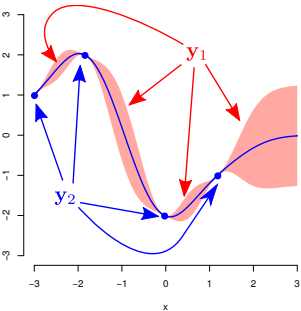
$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \Sigma_{\mathbf{f}\mathbf{f}^*} \\ \Sigma_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \Sigma_{\mathbf{f}\mathbf{f}^*} \\ \Sigma_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

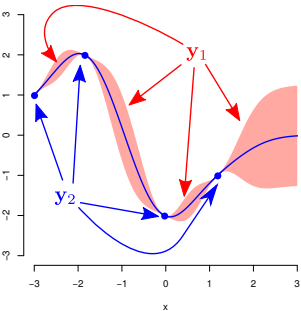


$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix} \right),$$

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \Sigma_{\mathbf{f}\mathbf{f}^*} \\ \Sigma_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



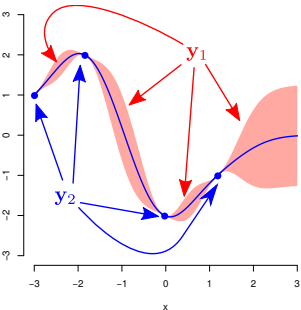
$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix} \right),$$

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N} \left(\mathbf{y}_1 \middle| \mathbf{a} + \mathbf{C}\mathbf{B}^{-1}(\mathbf{y}_2 - \mathbf{b}), \mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^T \right)$$

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{ff}} & \Sigma_{\mathbf{ff}^*} \\ \Sigma_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



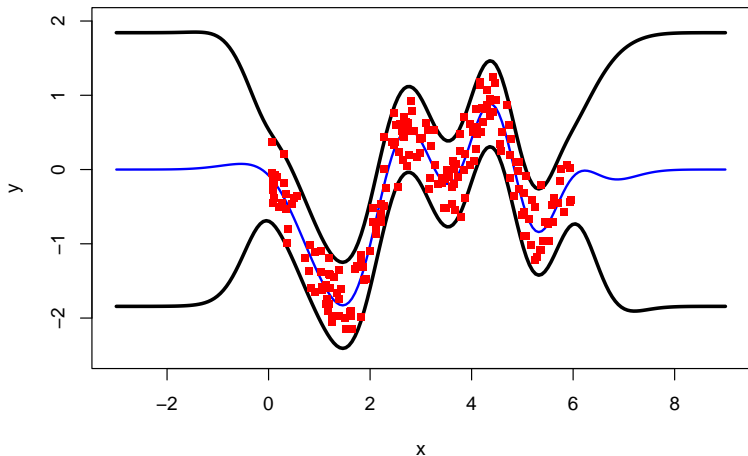
$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix} \right),$$

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N} \left(\mathbf{y}_1 \middle| \mathbf{a} + \mathbf{CB}^{-1}(\mathbf{y}_2 - \mathbf{b}), \mathbf{A} - \mathbf{CB}^{-1}\mathbf{C}^T \right)$$

$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \middle| \Sigma_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \Sigma_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \Sigma_{\mathbf{f}^*\mathbf{f}}^T \right)$$

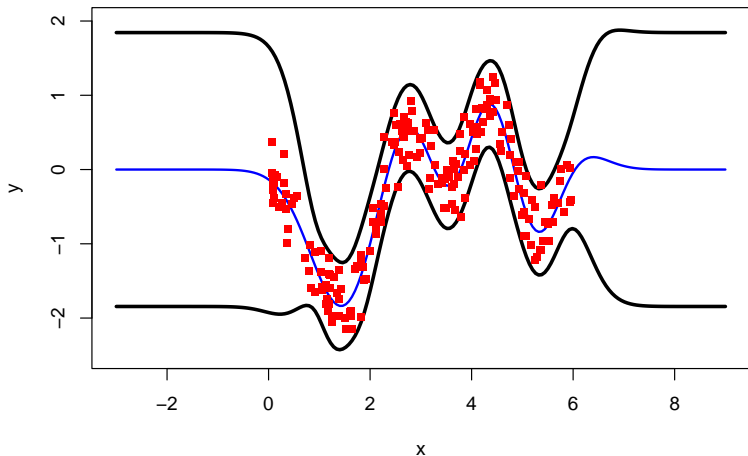
Nyström Approximation: Illustrative Example

Full GP



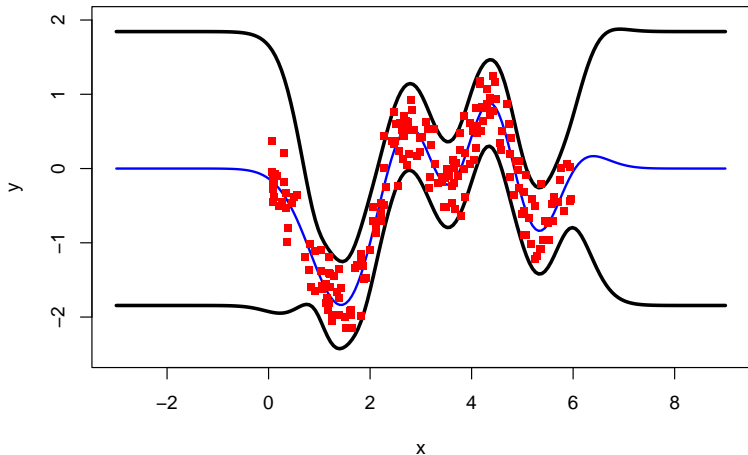
Nystrom Approximation: Illustrative Example

Nystrom GP ($M = 10$)



Nystrom Approximation: Illustrative Example

Nystrom GP ($M = 10$)



The approximation is similar to the full GP in some regions!

Summary of Nyström Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.

Summary of Nyström Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- If $M = N$ the method is exact since $\tilde{\Sigma} = \Sigma$.

Summary of Nyström Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- If $M = N$ the method is exact since $\tilde{\Sigma} = \Sigma$.
- For small M it can give bad results according to empirical evidence.

Summary of Nyström Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- If $M = N$ the method is exact since $\tilde{\Sigma} = \Sigma$.
- For small M it can give bad results according to empirical evidence.
- It can perform well if Σ is dominated by a few eigenvalues.

Summary of Nyström Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- If $M = N$ the method is exact since $\tilde{\Sigma} = \Sigma$.
- For small M it can give bad results according to empirical evidence.
- It can perform well if Σ is dominated by a few eigenvalues.
- As the M points are chosen at random it may give different results.

Summary of Nyström Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- If $M = N$ the method is exact since $\tilde{\Sigma} = \Sigma$.
- For small M it can give bad results according to empirical evidence.
- It can perform well if Σ is dominated by a few eigenvalues.
- As the M points are chosen at random it may give different results.
- Since the approximation is done only over the covariance matrix of the training data, negative predictive variances are possible, but rare.

Run the cells of the notebook to fit a sparse GP using the Nyström approximation and complete task 2!

Run the cells of the notebook to compare the exact and approximate predictive distribution and complete task 3!

Random Features Approximations

They can be used to approximate any stationary covariance function (only depends on the distance between points).

Random Features Approximations

They can be used to approximate any stationary covariance function (only depends on the distance between points).

Bochner's theorem:

A covariance function $C(\mathbf{x}, \mathbf{x}') = C(\mathbf{x} - \mathbf{x}')$ on \mathbb{R}^D is positive definite if and only if $C(\mathbf{x} - \mathbf{x}')$ is the Fourier transform of a distribution $s(\mathbf{w})$.

Random Features Approximations

They can be used to approximate any stationary covariance function (only depends on the distance between points).

Bochner's theorem:

A covariance function $C(\mathbf{x}, \mathbf{x}') = C(\mathbf{x} - \mathbf{x}')$ on \mathbb{R}^D is positive definite if and only if $C(\mathbf{x} - \mathbf{x}')$ is the Fourier transform of a distribution $s(\mathbf{w})$.

$$C(\mathbf{x}, \mathbf{x}') = \int \exp\{-i\mathbf{w}^\top(\mathbf{x} - \mathbf{x}')\} s(\mathbf{w}) d\mathbf{w},$$
$$s(\mathbf{w}) = \frac{1}{(2\pi)^D} \int \exp\{i\mathbf{w}^\top \boldsymbol{\tau}\} C(\boldsymbol{\tau}, \mathbf{0}) d\boldsymbol{\tau}.$$

Random Features Approximations

They can be used to approximate any stationary covariance function (only depends on the distance between points).

Bochner's theorem:

A covariance function $C(\mathbf{x}, \mathbf{x}') = C(\mathbf{x} - \mathbf{x}')$ on \mathbb{R}^D is positive definite if and only if $C(\mathbf{x} - \mathbf{x}')$ is the Fourier transform of a distribution $s(\mathbf{w})$.

$$C(\mathbf{x}, \mathbf{x}') = \int \exp\{-i\mathbf{w}^\top(\mathbf{x} - \mathbf{x}')\} s(\mathbf{w}) d\mathbf{w},$$
$$s(\mathbf{w}) = \frac{1}{(2\pi)^D} \int \exp\{i\mathbf{w}^\top \boldsymbol{\tau}\} C(\boldsymbol{\tau}, \mathbf{0}) d\boldsymbol{\tau}.$$

$s(\mathbf{w})$ is called the spectral density of the covariance function.

Covariances as Expectations of Cosines

Due to Bochner's theorem, the covariance can be written as:

$$\begin{aligned} C(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{s(\mathbf{w})} \left[\exp\{-i\mathbf{w}^T(\mathbf{x} - \mathbf{x}')\} \right] \\ &= 2\mathbb{E}_{s(\mathbf{w}), b \sim U[0, 2\pi]} \left[\cos(\mathbf{w}^T \mathbf{x} + b) \cos(\mathbf{w}^T \mathbf{x}' + b) \right]. \end{aligned}$$

Covariances as Expectations of Cosines

Due to Bochner's theorem, the covariance can be written as:

$$\begin{aligned} C(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{s(\mathbf{w})} \left[\exp\{-i\mathbf{w}^T(\mathbf{x} - \mathbf{x}')\} \right] \\ &= 2\mathbb{E}_{s(\mathbf{w}), b \sim U[0, 2\pi]} \left[\cos(\mathbf{w}^T \mathbf{x} + b) \cos(\mathbf{w}^T \mathbf{x}' + b) \right]. \end{aligned}$$

The expectation can be approximated by a Monte Carlo average!

Covariances as Expectations of Cosines

Due to Bochner's theorem, the covariance can be written as:

$$\begin{aligned} C(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{s(\mathbf{w})} \left[\exp\{-i\mathbf{w}^T(\mathbf{x} - \mathbf{x}')\} \right] \\ &= 2\mathbb{E}_{s(\mathbf{w}), b \sim U[0, 2\pi]} \left[\cos(\mathbf{w}^T \mathbf{x} + b) \cos(\mathbf{w}^T \mathbf{x}' + b) \right]. \end{aligned}$$

The expectation can be approximated by a Monte Carlo average!

We can reduce the variance of the estimator by generating M samples:

$$C(\mathbf{x}, \mathbf{x}') \approx \frac{2}{M} \sum_{m=1}^M \cos(\mathbf{w}_m^T \mathbf{x} + b_m) \cos(\mathbf{w}_m^T \mathbf{x}' + b_m) = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

with $\phi(\mathbf{x}) = \sqrt{\frac{2}{M}} \cos(\mathbf{W}^T \mathbf{x} + \mathbf{b})$ a random M feature expansion.

Covariances as Expectations of Cosines

Due to Bochner's theorem, the covariance can be written as:

$$\begin{aligned} C(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{s(\mathbf{w})} \left[\exp\{-i\mathbf{w}^T(\mathbf{x} - \mathbf{x}')\} \right] \\ &= 2\mathbb{E}_{s(\mathbf{w}), b \sim U[0, 2\pi]} \left[\cos(\mathbf{w}^T \mathbf{x} + b) \cos(\mathbf{w}^T \mathbf{x}' + b) \right]. \end{aligned}$$

The expectation can be approximated by a Monte Carlo average!

We can reduce the variance of the estimator by generating M samples:

$$C(\mathbf{x}, \mathbf{x}') \approx \frac{2}{M} \sum_{m=1}^M \cos(\mathbf{w}_m^T \mathbf{x} + b_m) \cos(\mathbf{w}_m^T \mathbf{x}' + b_m) = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

with $\phi(\mathbf{x}) = \sqrt{\frac{2}{M}} \cos(\mathbf{W}^T \mathbf{x} + \mathbf{b})$ a random M feature expansion.

For the squared exponential covariance function $s(\mathbf{w})$ is Gaussian!

Approximate Covariance Function

The covariance matrix can be simply approximated as:

$$\Sigma \approx \tilde{\Sigma} = \Phi\Phi^T$$

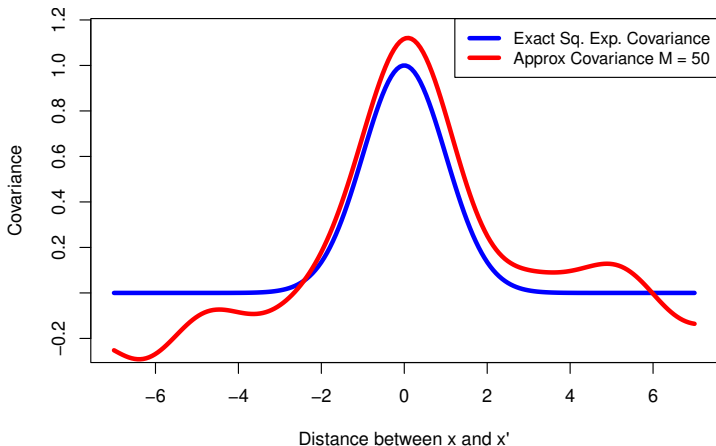
and hence $\mathbf{I}\sigma^2 + \tilde{\Sigma}$ can be inverted with cost $\mathcal{O}(NM^2)$.

Approximate Covariance Function

The covariance matrix can be simply approximated as:

$$\Sigma \approx \tilde{\Sigma} = \Phi\Phi^T$$

and hence $\mathbf{I}\sigma^2 + \tilde{\Sigma}$ can be inverted with cost $\mathcal{O}(NM^2)$.



Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}\mathbf{f}^*} \\ \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}\mathbf{f}^*} \\ \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

All prior covariances are now approximated using dot products with the random features computed before!

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{ff}} & \tilde{\Sigma}_{\mathbf{ff}^*} \\ \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

All prior covariances are now approximated using dot products with the random features computed before!

$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{f}, \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}^*} - \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}}^T \right)$$

Predictive Distribution

We want to compute the value of f^* at a new \mathbf{x}^* :

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}\mathbf{f}^*} \\ \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} & \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

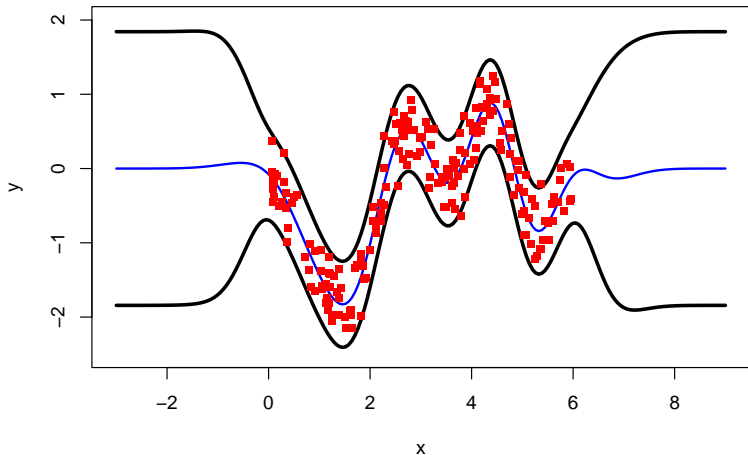
All prior covariances are now approximated using dot products with the random features computed before!

$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{f}, \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}^*} - \tilde{\Sigma}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}^*}^T \right)$$

The computational cost is $\mathcal{O}(NM^2)$!

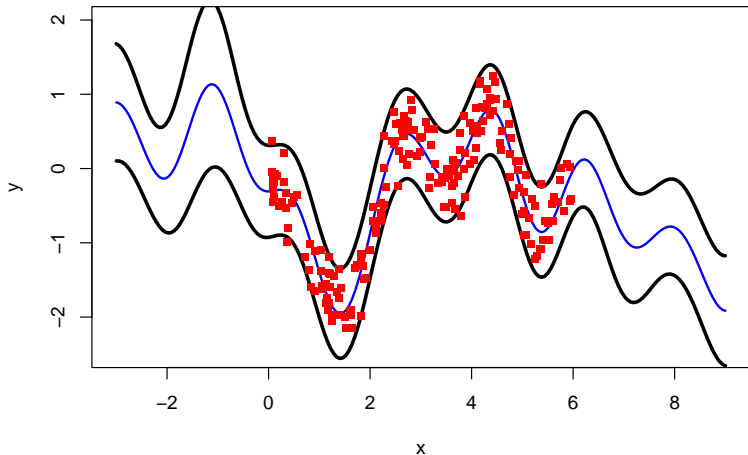
Random Features: Illustrative Example

Full GP



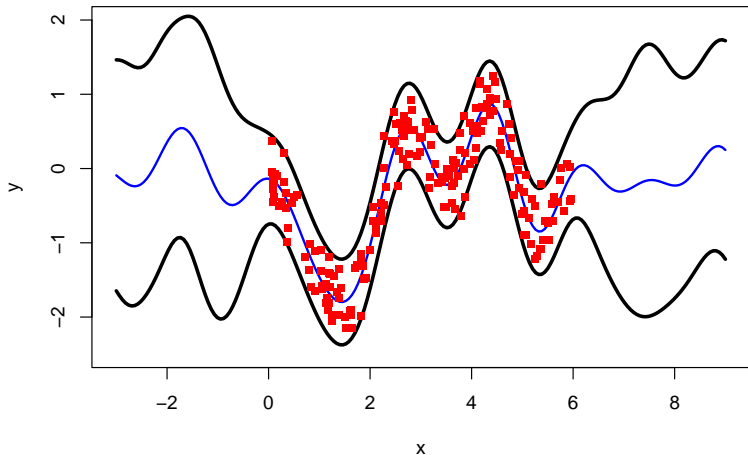
Random Features: Illustrative Example

Random Features GP ($M = 10$)



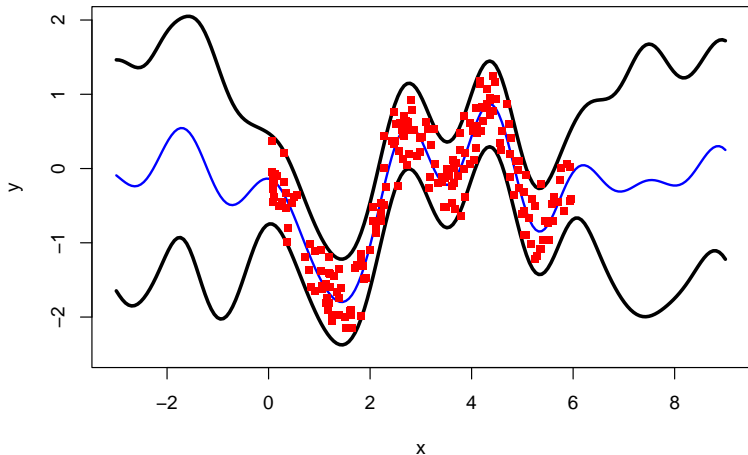
Random Features: Illustrative Example

Random Features GP ($M = 50$)



Random Features: Illustrative Example

Random Features GP ($M = 50$)



In regions with no data the approximation may be wiggling a lot!

Summary of Random Features Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.

Summary of Random Features Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- For small M it can give bad results due to the wiggling effect of cosine features.

Summary of Random Features Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- For small M it can give bad results due to the wiggling effect of cosine features.
- Guaranteed to be exact only for $M \rightarrow \infty$.

Summary of Random Features Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- For small M it can give bad results due to the wiggling effect of cosine features.
- Guaranteed to be exact only for $M \rightarrow \infty$.
- It is restricted to stationary covariance functions.

Summary of Random Features Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- For small M it can give bad results due to the wiggling effect of cosine features.
- Guaranteed to be exact only for $M \rightarrow \infty$.
- It is restricted to stationary covariance functions.
- Very simple to implement!

Summary of Random Features Approximation

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- For small M it can give bad results due to the wiggling effect of cosine features.
- Guaranteed to be exact only for $M \rightarrow \infty$.
- It is restricted to stationary covariance functions.
- Very simple to implement!
- Equivalent to a neural network with a hidden layer with M units and cosine activations, and a Bayesian linear model in the last layer!

**Run the cells of the notebook to fit a sparse GP
using the Random Features approximation and
complete task 4!**

Full Independent Training Conditional (FITC)

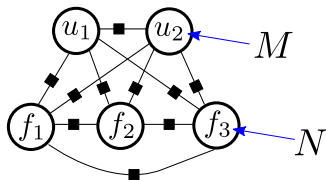
1. Extend model with $M \ll N$ inducing points and outputs at $\bar{\mathbf{X}}$.

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{ff}} & \Sigma_{\mathbf{fu}} \\ \Sigma_{\mathbf{uf}} & \Sigma_{\mathbf{uu}} \end{bmatrix} \right)$$

Full Independent Training Conditional (FITC)

1. Extend model with $M \ll N$ inducing points and outputs at $\bar{\mathbf{X}}$.

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{ff} & \Sigma_{fu} \\ \Sigma_{uf} & \Sigma_{uu} \end{bmatrix} \right)$$



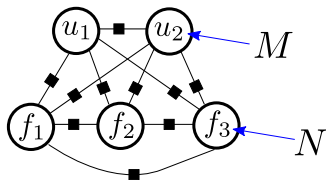
Full Independent Training Conditional (FITC)

1. Extend model with $M \ll N$ inducing points and outputs at $\bar{\mathbf{X}}$.

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{ff}} & \Sigma_{\mathbf{fu}} \\ \Sigma_{\mathbf{uf}} & \Sigma_{\mathbf{uu}} \end{bmatrix} \right)$$

2. Introduce conditional independences:

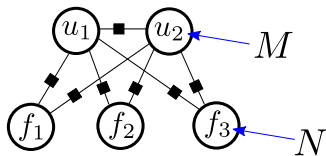
$$p(\mathbf{f}|\mathbf{u}) = \prod_{i=1}^N p(f_i|\mathbf{u})$$



Full Independent Training Conditional (FITC)

1. Extend model with $M \ll N$ inducing points and outputs at $\bar{\mathbf{X}}$.

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{ff}} & \Sigma_{\mathbf{fu}} \\ \Sigma_{\mathbf{uf}} & \Sigma_{\mathbf{uu}} \end{bmatrix} \right)$$



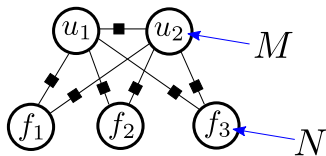
2. Introduce conditional independences:

$$p(\mathbf{f}|\mathbf{u}) = \prod_{i=1}^N p(f_i|\mathbf{u})$$

Full Independent Training Conditional (FITC)

1. Extend model with $M \ll N$ inducing points and outputs at $\bar{\mathbf{X}}$.

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{ff}} & \Sigma_{\mathbf{fu}} \\ \Sigma_{\mathbf{uf}} & \Sigma_{\mathbf{uu}} \end{bmatrix} \right) \right)$$



2. Introduce conditional independences:

$$p(\mathbf{f}|\mathbf{u}) = \prod_{i=1}^N p(f_i|\mathbf{u})$$

3. Marginalize \mathbf{u} to obtain an approximate GP prior for \mathbf{f} .

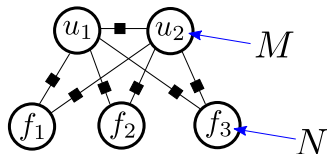
$$p(\mathbf{f}) = \int p(\mathbf{f}|\mathbf{u})p(\mathbf{u})d\mathbf{u} = \prod_{i=1}^N p(f_i|\mathbf{u})p(\mathbf{u})d\mathbf{u} = \mathcal{N}(\mathbf{f}|\mathbf{0}, \tilde{\Sigma}_{\mathbf{ff}})$$

where $\tilde{\Sigma}_{\mathbf{ff}} = \mathbf{D} + \mathbf{Q}_{\mathbf{ff}}$ with \mathbf{D} diagonal and $\mathbf{Q}_{\mathbf{ff}} = \Sigma_{\mathbf{fu}}\Sigma_{\mathbf{uu}}^{-1}\Sigma_{\mathbf{uf}}$ of rank M .

Full Independent Training Conditional (FITC)

4. We make the prediction of f^* at \mathbf{x}^* by considering the approximate GP prior:

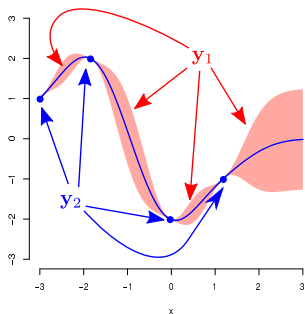
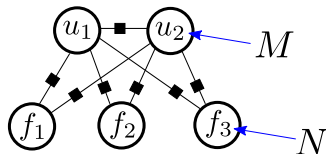
$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} \\ \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



Full Independent Training Conditional (FITC)

4. We make the prediction of f^* at \mathbf{x}^* by considering the approximate GP prior:

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} \\ \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



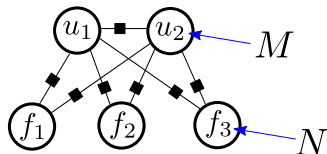
$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix} \right),$$

$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)},$$

Full Independent Training Conditional (FITC)

4. We make the prediction of f^* at \mathbf{x}^* by considering the approximate GP prior:

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} & \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} \\ \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$

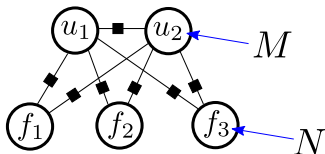


$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{Q}_{\mathbf{f}^*\mathbf{f}}^T \right)$$

Full Independent Training Conditional (FITC)

4. We make the prediction of f^* at \mathbf{x}^* by considering the approximate GP prior:

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{ff}} & \mathbf{Q}_{\mathbf{ff}^*} \\ \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



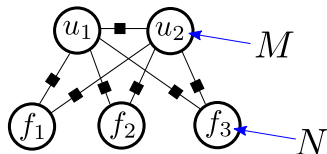
$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{Q}_{\mathbf{f}^*\mathbf{f}}^T \right)$$

Due to the structure in $\tilde{\Sigma}_{\mathbf{ff}}$ all computations have cost in $\mathcal{O}(NM^2)$.

Full Independent Training Conditional (FITC)

4. We make the prediction of f^* at \mathbf{x}^* by considering the approximate GP prior:

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{ff}} & \mathbf{Q}_{\mathbf{ff}^*} \\ \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{Q}_{\mathbf{f}^*\mathbf{f}}^T \right)$$

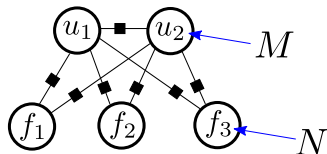
Due to the structure in $\tilde{\Sigma}_{\mathbf{ff}}$ all computations have cost in $\mathcal{O}(NM^2)$.

5. How do we find the location of the inducing points $\bar{\mathbf{X}}$?

Full Independent Training Conditional (FITC)

4. We make the prediction of f^* at \mathbf{x}^* by considering the approximate GP prior:

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{\mathbf{ff}} & \mathbf{Q}_{\mathbf{ff}^*} \\ \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} & \Sigma_{\mathbf{f}^*\mathbf{f}^*} \end{bmatrix} \right)$$



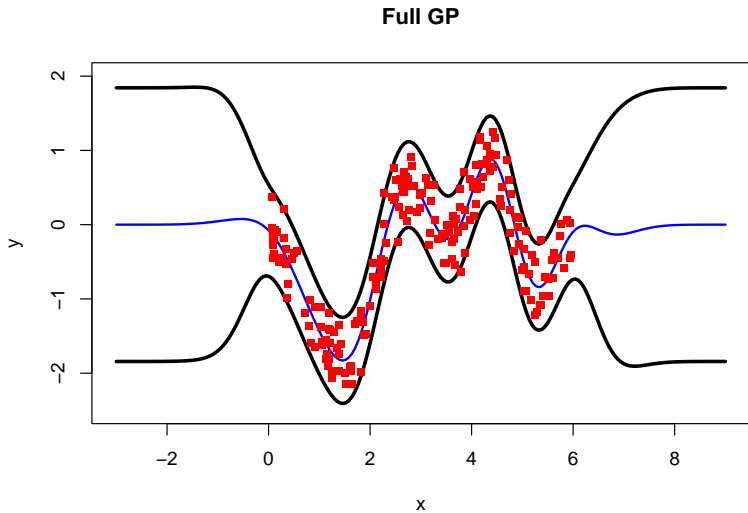
$$p(\mathbf{f}^* | \mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{ff}}^{-1} \mathbf{Q}_{\mathbf{f}^*\mathbf{f}}^T \right)$$

Due to the structure in $\tilde{\Sigma}_{\mathbf{ff}}$ all computations have cost in $\mathcal{O}(NM^2)$.

5. How do we find the location of the inducing points $\bar{\mathbf{X}}$?

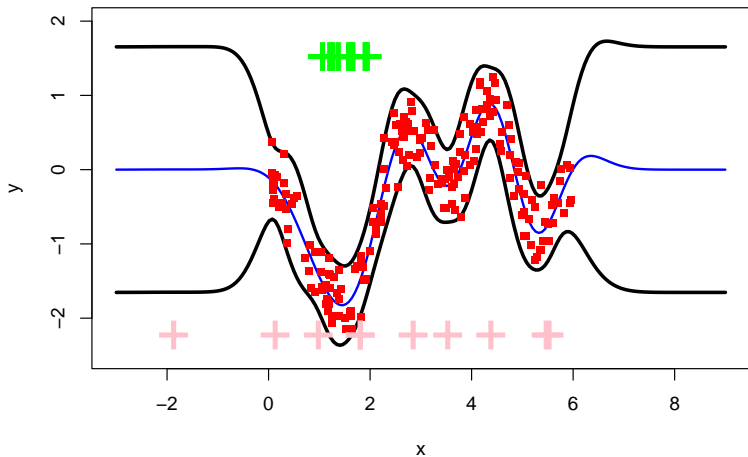
Simply treat them as prior parameters and maximize the approximate marginal likelihood $p(\mathbf{f} | \mathbf{0}, \tilde{\Sigma}_{\mathbf{ff}})$!

FITC: Illustrative Example



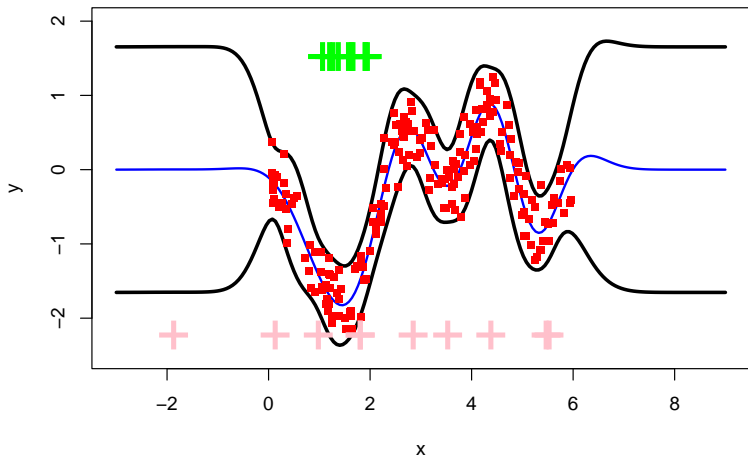
FITC: Illustrative Example

FITC ($M = 10$)



FITC: Illustrative Example

FITC ($M = 10$)



The inducing points cover the regions where the function changes!

Summary of FITC

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.

Summary of FITC

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.

Summary of FITC

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.
- Guaranteed to be exact if $M = N$ and the inducing points are not optimized and located at the training points.

Summary of FITC

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.
- Guaranteed to be exact if $M = N$ and the inducing points are not optimized and located at the training points.
- It can be understood as considering heteroscedastic (input dependent) noise!

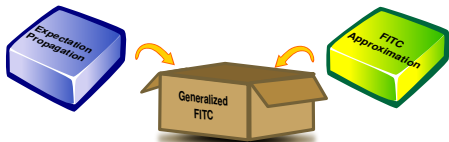
**Run the cells of the notebook to fit a sparse GP
using the FITC approximation and complete task
5!**

Generalized FITC

Combines FITC with the use of Expectation Propagation to address binary classification problems!

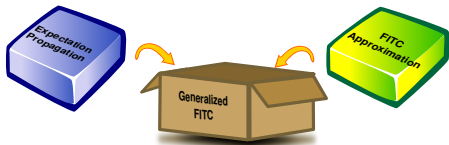
Generalized FITC

Combines FITC with the use of Expectation Propagation to address binary classification problems!



Generalized FITC

Combines FITC with the use of Expectation Propagation to address binary classification problems!

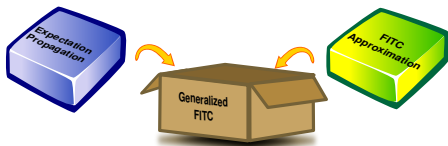


Assumes $y_i \in \{-1, 1\}$ and a probit likelihood:

$$p(y_i | f(\mathbf{x}_i)) = \phi(y_i f(\mathbf{x}_i)), \quad \phi(\cdot) \equiv \text{The c.d.f. of a standard Gaussian.}$$

Generalized FITC

Combines FITC with the use of Expectation Propagation to address binary classification problems!



Assumes $y_i \in \{-1, 1\}$ and a probit likelihood:

$$p(y_i | f(\mathbf{x}_i)) = \phi(y_i f(\mathbf{x}_i)), \quad \phi(\cdot) \equiv \text{The c.d.f. of a standard Gaussian.}$$

Approximates with a **Gaussian distribution** the intractable posterior:

$$p(\mathbf{f} | \mathbf{y}) = \frac{\prod_{i=1}^N \phi(y_i f(\mathbf{x}_i)) \mathcal{N}(\mathbf{f} | \mathbf{0}, \tilde{\Sigma})}{p(\mathbf{y})},$$

where $\tilde{\Sigma}$ is the **approximate FITC covariance matrix**.

Introduction to Expectation Propagation

Approximates an intractable distribution p by a parametric distribution q .

Introduction to Expectation Propagation

Approximates an intractable distribution p by a parametric distribution q .

It is based on the minimization of the KL-divergence, $\text{KL}(p||q)$:

$$\int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} = \text{KL}(q|p) \geq 0.$$

Introduction to Expectation Propagation

Approximates an intractable distribution p by a parametric distribution q .

It is based on the minimization of the KL-divergence, $\text{KL}(p||q)$:

$$\int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} = \text{KL}(q|p) \geq 0.$$

q is restricted to belong to a family of distributions **closed under the product and ratio operation**: The **exponential family**.

Introduction to Expectation Propagation

Approximates an intractable distribution p by a parametric distribution q .

It is based on the minimization of the KL-divergence, $\text{KL}(p||q)$:

$$\int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} = \text{KL}(q|p) \geq 0.$$

q is restricted to belong to a family of distributions **closed under the product and ratio operation**: The **exponential family**.

The exponential family:

$$q(\mathbf{x}) = \exp \left(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) - g(\boldsymbol{\eta}) \right), \quad g(\boldsymbol{\eta}) = \log \int \exp \left(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right) d\mathbf{x}$$

where $\boldsymbol{\eta}$ is a vector of natural parameters of q , $\mathbf{u}(\mathbf{x})$ are the sufficient statistics and $g(\boldsymbol{\eta})$ is a **log partition function**.

Examples of Distributions in the Exponential Family

Gaussian:

$$\mathcal{N}(x|\mu, \sigma^2) = 1/\sqrt{2\pi\sigma^2} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

Examples of Distributions in the Exponential Family

Gaussian:

$$\mathcal{N}(x|\mu, \sigma^2) = 1/\sqrt{2\pi\sigma^2} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

Exponential form:

$$\mathcal{N}(x|\mu, \sigma^2) = \exp\left(\boldsymbol{\eta}^\top \mathbf{u}(x) - g(\boldsymbol{\eta})\right)$$

$$\boldsymbol{\eta} = (\mu/\sigma^2, 1.0/\sigma^2)^\top, \quad \mathbf{u}(x) = (x, -0.5x^2)^\top, \quad g(\boldsymbol{\eta}) = -\frac{1}{2} \log \frac{2\pi}{\eta_2} + \frac{\eta_1^2}{2\eta_2}.$$

Examples of Distributions in the Exponential Family

Gaussian:

$$\mathcal{N}(x|\mu, \sigma^2) = 1/\sqrt{2\pi\sigma^2} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

Exponential form:

$$\mathcal{N}(x|\mu, \sigma^2) = \exp\left(\boldsymbol{\eta}^\top \mathbf{u}(x) - g(\boldsymbol{\eta})\right)$$

$$\boldsymbol{\eta} = (\mu/\sigma^2, 1.0/\sigma^2)^\top, \quad \mathbf{u}(x) = (x, -0.5x^2)^\top, \quad g(\boldsymbol{\eta}) = -\frac{1}{2} \log \frac{2\pi}{\eta_2} + \frac{\eta_1^2}{2\eta_2}.$$

Most parametric distributions belong to the exponential family!

Product and Ratio of Gaussians

Consider these two Gaussian distributions:

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ -\frac{1}{2}\sigma_1^2(x - \mu_1)^2 \right\},$$

$$p_2(x) = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ -\frac{1}{2}\sigma_2^2(x - \mu_2)^2 \right\}.$$

Product and Ratio of Gaussians

Consider these two Gaussian distributions:

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ -\frac{1}{2}\sigma_1^2(x - \mu_1)^2 \right\},$$

$$p_2(x) = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ -\frac{1}{2}\sigma_2^2(x - \mu_2)^2 \right\}.$$

- $p_1(x)p_2(x)$ is Gaussian with natural parameters $\eta_1 + \eta_2$.

Product and Ratio of Gaussians

Consider these two Gaussian distributions:

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ -\frac{1}{2}\sigma_1^2(x - \mu_1)^2 \right\},$$

$$p_2(x) = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ -\frac{1}{2}\sigma_2^2(x - \mu_2)^2 \right\}.$$

- $p_1(x)p_2(x)$ is Gaussian with natural parameters $\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2$.
- The log-normalization constant of $p_1(x)p_2(x)$ is $g(\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2) - g(\boldsymbol{\eta}_1) - g(\boldsymbol{\eta}_2)$.

Product and Ratio of Gaussians

Consider these two Gaussian distributions:

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ -\frac{1}{2}\sigma_1^2(x - \mu_1)^2 \right\},$$

$$p_2(x) = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ -\frac{1}{2}\sigma_2^2(x - \mu_2)^2 \right\}.$$

- $p_1(x)p_2(x)$ is Gaussian with natural parameters $\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2$.
- The log-normalization constant of $p_1(x)p_2(x)$ is $g(\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2) - g(\boldsymbol{\eta}_1) - g(\boldsymbol{\eta}_2)$.
- $p_1(x)/p_2(x)$ is Gaussian with natural parameters $\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2$.

Product and Ratio of Gaussians

Consider these two Gaussian distributions:

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ -\frac{1}{2}\sigma_1^2(x - \mu_1)^2 \right\},$$

$$p_2(x) = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ -\frac{1}{2}\sigma_2^2(x - \mu_2)^2 \right\}.$$

- $p_1(x)p_2(x)$ is Gaussian with natural parameters $\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2$.
- The log-normalization constant of $p_1(x)p_2(x)$ is $g(\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2) - g(\boldsymbol{\eta}_1) - g(\boldsymbol{\eta}_2)$.
- $p_1(x)/p_2(x)$ is Gaussian with natural parameters $\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2$.
- The log-normalization constant of $p_1(x)/p_2(x)$ is $g(\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2) - g(\boldsymbol{\eta}_1) + g(\boldsymbol{\eta}_2)$.

KL-Divergence Minimization

Consider the KL-divergence between p and q (q in the exponential family):

$$\text{KL}(p||q) = - \int p(\mathbf{x}) \log \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} = g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_p[\mathbf{u}(\mathbf{x})] + \text{Const} .$$

KL-Divergence Minimization

Consider the KL-divergence between p and q (q in the exponential family):

$$\text{KL}(p||q) = - \int p(\mathbf{x}) \log \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} = g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_p[\mathbf{u}(\mathbf{x})] + \text{Const} .$$

When **minimizing** with respect to the natural parameters $\boldsymbol{\eta}$ of q :

$$\frac{\partial \text{KL}(p||q)}{\partial \boldsymbol{\eta}} = 0 \iff \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_p[\mathbf{u}(\mathbf{x})] ,$$

KL-Divergence Minimization

Consider the KL-divergence between p and q (q in the exponential family):

$$\text{KL}(p||q) = - \int p(\mathbf{x}) \log \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} = g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_p[\mathbf{u}(\mathbf{x})] + \text{Const} .$$

When **minimizing** with respect to the natural parameters $\boldsymbol{\eta}$ of q :

$$\frac{\partial \text{KL}(p||q)}{\partial \boldsymbol{\eta}} = 0 \iff \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_p[\mathbf{u}(\mathbf{x})] ,$$

Furthermore, it is possible to show that:

$$\frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_q[\mathbf{u}(\mathbf{x})] .$$

KL-Divergence Minimization

Consider the KL-divergence between p and q (q in the exponential family):

$$\text{KL}(p||q) = - \int p(\mathbf{x}) \log \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} = g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_p[\mathbf{u}(\mathbf{x})] + \text{Const.}$$

When **minimizing** with respect to the natural parameters $\boldsymbol{\eta}$ of q :

$$\frac{\partial \text{KL}(p||q)}{\partial \boldsymbol{\eta}} = 0 \iff \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_p[\mathbf{u}(\mathbf{x})],$$

Furthermore, it is possible to show that:

$$\frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_q[\mathbf{u}(\mathbf{x})].$$

KL($p||q$) is minimized by matching expected sufficient statistics.

KL-Divergence Minimization

Consider the KL-divergence between p and q (q in the exponential family):

$$\text{KL}(p||q) = - \int p(\mathbf{x}) \log \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} = g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_p[\mathbf{u}(\mathbf{x})] + \text{Const.}$$

When **minimizing** with respect to the natural parameters $\boldsymbol{\eta}$ of q :

$$\frac{\partial \text{KL}(p||q)}{\partial \boldsymbol{\eta}} = 0 \iff \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_p[\mathbf{u}(\mathbf{x})],$$

Furthermore, it is possible to show that:

$$\frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \mathbb{E}_q[\mathbf{u}(\mathbf{x})].$$

KL($p||q$) is minimized by matching expected sufficient statistics.

If q is Gaussian, then we have to match $\mathbb{E}_q[\mathbf{x}] = \mathbb{E}_p[\mathbf{x}]$ and $\mathbb{E}_q[\mathbf{x}\mathbf{x}^T] = \mathbb{E}_p[\mathbf{x}\mathbf{x}^T]$.

Joint Approximation

EP approximates this joint distribution by a **product of simpler factors**:

$$p(\mathbf{f}, \mathbf{y}) = \prod_{i=1}^N \phi_i(y_i f(\mathbf{x}_i) \mathcal{N}(\mathbf{f} | \mathbf{0}, \tilde{\Sigma})) = \prod_i t_i(\mathbf{f}) \approx \prod_i \tilde{t}_i(\mathbf{f}),$$

where each \tilde{t}_i approximates the corresponding t_i . Each \tilde{t}_i must **belong to the exponential family** but **need not be normalized**.

Joint Approximation

EP approximates this joint distribution by a **product of simpler factors**:

$$p(\mathbf{f}, \mathbf{y}) = \prod_{i=1}^N \phi_i(y_i f(\mathbf{x}_i) | \mathcal{N}(\mathbf{f} | \mathbf{0}, \tilde{\Sigma})) = \prod_i t_i(\mathbf{f}) \approx \prod_i \tilde{t}_i(\mathbf{f}),$$

where each \tilde{t}_i approximates the corresponding t_i . Each \tilde{t}_i must **belong to the exponential family** but **need not be normalized**.

The exponential family is closed under the product and $\prod_i \tilde{t}_i$ can be easily normalized to compute an approximate distribution:

$$p(\mathbf{f} | \mathbf{y}) = \frac{1}{p(\mathbf{y})} \prod_i t_i(\mathbf{f}) \approx \frac{1}{Z} \prod_i \tilde{t}_i(\mathbf{f}) = q(\mathbf{f}),$$

where $Z = \int \prod_i \tilde{t}_i(\mathbf{f}) d\mathbf{f}$ can be used to **approximate** $p(\mathbf{y})$.

Joint Approximation

EP approximates this joint distribution by a **product of simpler factors**:

$$p(\mathbf{f}, \mathbf{y}) = \prod_{i=1}^N \phi_i(y_i f(\mathbf{x}_i) | \mathcal{N}(\mathbf{f} | \mathbf{0}, \tilde{\Sigma})) = \prod_i t_i(\mathbf{f}) \approx \prod_i \tilde{t}_i(\mathbf{f}),$$

where each \tilde{t}_i approximates the corresponding t_i . Each \tilde{t}_i must **belong to the exponential family** but **need not be normalized**.

The exponential family is closed under the product and $\prod_i \tilde{t}_i$ can be easily normalized to compute an approximate distribution:

$$p(\mathbf{f} | \mathbf{y}) = \frac{1}{p(\mathbf{y})} \prod_i t_i(\mathbf{f}) \approx \frac{1}{Z} \prod_i \tilde{t}_i(\mathbf{f}) = q(\mathbf{f}),$$

where $Z = \int \prod_i \tilde{t}_i(\mathbf{f}) d\mathbf{f}$ can be used to **approximate** $p(\mathbf{y})$.

Therefore q has the same form as the approximate factors!

Approximate Factors

How do we determine each approximate factor \tilde{t}_i ?

Approximate Factors

How do we determine each approximate factor \tilde{t}_i ?

We would like to **minimize** $\text{KL}(p||q)$, but this is **intractable!**

Approximate Factors

How do we determine each approximate factor \tilde{t}_i ?

We would like to **minimize** $\text{KL}(p||q)$, but this is **intractable**!

EP minimizes the KL divergence **between pairs** of t_i and \tilde{t}_i . This has the risk that **the product** may not be a good approximation. EP tries to **circumvent** this by an iterative procedure.

Approximate Factors

How do we determine each approximate factor \tilde{t}_i ?

We would like to **minimize** $\text{KL}(p||q)$, but this is **intractable**!

EP minimizes the KL divergence **between pairs** of t_i and \tilde{t}_i . This has the risk that **the product** may not be a good approximation. EP tries to **circumvent** this by an iterative procedure.

Suppose we wish to refine \tilde{t}_j . We first remove this factor from the product:

$$q^{\setminus j}(\mathbf{f}) \propto \prod_{i \neq j} \tilde{t}_i(\mathbf{f}) \propto q(\mathbf{f}) / \tilde{t}_j(\mathbf{f}),$$

Approximate Factors

How do we determine each approximate factor \tilde{t}_i ?

We would like to **minimize** $\text{KL}(p||q)$, but this is **intractable!**

EP minimizes the KL divergence **between pairs** of t_i and \tilde{t}_i . This has the risk that **the product** may not be a good approximation. EP tries to **circumvent** this by an iterative procedure.

Suppose we wish to refine \tilde{t}_j . We first remove this factor from the product:

$$q^{\setminus j}(\mathbf{f}) \propto \prod_{i \neq j} \tilde{t}_i(\mathbf{f}) \propto q(\mathbf{f}) / \tilde{t}_j(\mathbf{f}),$$

Then, \tilde{t}_j is updated to minimize the KL-divergence between:

$$q_{\text{new}}(\mathbf{f}) \propto \tilde{t}_j(\mathbf{f}) q^{\setminus j}(\mathbf{f}), \quad \hat{p}_j(\mathbf{f}) = \frac{1}{Z_j} t_j(\mathbf{f}) q^{\setminus j}(\mathbf{f}), \quad Z_j = \int t_j(\mathbf{f}) q^{\setminus j}(\mathbf{f}) d\mathbf{f},$$

where $q^{\setminus j}$ is fixed. This ensures that \tilde{t}_j is accurate where $q^{\setminus j}$ is high.

Approximate Factors

In practice, \tilde{t}_j is found by first **minimizing** with respect to q_{new} :

$$\text{KL} \left(\frac{t_j(\mathbf{f})q^{\setminus j}(\mathbf{f})}{Z_j} \middle| q_{\text{new}}(\mathbf{f}) \right) .$$

Approximate Factors

In practice, \tilde{t}_j is found by first **minimizing** with respect to q_{new} :

$$\text{KL} \left(\frac{t_j(\mathbf{f})q^{\setminus j}(\mathbf{f})}{Z_j} \middle| q_{\text{new}}(\mathbf{f}) \right) .$$

This is done by **matching expected sufficient statistics**. As q is Gaussian, we only have to match the mean and the variance.

Approximate Factors

In practice, \tilde{t}_j is found by first **minimizing** with respect to q_{new} :

$$\text{KL} \left(\frac{t_j(\mathbf{f})q^{\setminus j}(\mathbf{f})}{Z_j} \middle| q_{\text{new}}(\mathbf{f}) \right) .$$

This is done by **matching expected sufficient statistics**. As q is Gaussian, we only have to match the mean and the variance.

It is required that the moments of $\hat{p}_j(\mathbf{f}) = 1/Z_j f_j(\mathbf{f})q^{\setminus j}(\mathbf{f})$ are tractable.

Approximate Factors

In practice, \tilde{t}_j is found by first **minimizing** with respect to q_{new} :

$$\text{KL} \left(\frac{t_j(\mathbf{f})q^{\vee}(\mathbf{f})}{Z_j} \middle| q_{\text{new}}(\mathbf{f}) \right).$$

This is done by **matching expected sufficient statistics**. As q is Gaussian, we only have to match the mean and the variance.

It is required that the moments of $\hat{p}_j(\mathbf{f}) = 1/Z_j t_j(\mathbf{f})q^{\vee}(\mathbf{f})$ are tractable.

The refined factor \tilde{t}_j is set in practice to be:

$$\tilde{t}_j(\mathbf{f}) = Z_j \frac{q_{\text{new}}(\mathbf{f})}{q^{\vee}(\mathbf{f})}, \quad \text{with} \quad \tilde{t}_j(\mathbf{f})q^{\vee}(\mathbf{f}) \propto q_{\text{new}},$$

which ensures that $\tilde{t}_j(\mathbf{f})q^{\vee}(\mathbf{f})$ and $t_j(\mathbf{f})q^{\vee}(\mathbf{f})$ **integrate the same**.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :
 - 1 Choose a factor \tilde{t}_j to refine.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :
 - 1 Choose a factor \tilde{t}_j to refine.
 - 2 Remove \tilde{t}_j from q by division $q^{\setminus j} \propto q / \tilde{t}_j$.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :
 - 1 Choose a factor \tilde{t}_j to refine.
 - 2 Remove \tilde{t}_j from q by division $q^{\setminus j} \propto q / \tilde{t}_j$.
 - 3 Compute Z_j and \hat{p}_j and find q_{new} by minimizing $\text{KL}(\hat{p}_j || q_{\text{new}})$.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :
 - 1 Choose a factor \tilde{t}_j to refine.
 - 2 Remove \tilde{t}_j from q by division $q^{\setminus j} \propto q / \tilde{t}_j$.
 - 3 Compute Z_j and \hat{p}_j and find q_{new} by minimizing $\text{KL}(\hat{p}_j || q_{\text{new}})$.
 - 4 Compute and store the new factor $\tilde{t}_j = Z_j q_{\text{new}} / q^{\setminus j}$.

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :
 - 1 Choose a factor \tilde{t}_j to refine.
 - 2 Remove \tilde{t}_j from q by division $q^{\setminus j} \propto q/\tilde{t}_j$.
 - 3 Compute Z_j and \hat{p}_j and find q_{new} by minimizing $\text{KL}(\hat{p}_j||q_{\text{new}})$.
 - 4 Compute and store the new factor $\tilde{t}_j = Z_j q_{\text{new}}/q^{\setminus j}$.
- 3 Evaluate the approximation to the model evidence:

$$p(\mathbf{y}) \approx Z = \int \prod_j \tilde{t}_j(\mathbf{f}) d\mathbf{f}.$$

Full Algorithm

Several passes are made through the factors until they **converge**. The **model evidence** is approximated by the **normalizing constant** of q .

EP Algorithm in General: Computes q and an approximation to $p(\mathbf{y})$.

- 1 Initialize q and each \tilde{t}_i to be uniform.
- 2 Repeat until convergence of the \tilde{t}_i :
 - 1 Choose a factor \tilde{t}_j to refine.
 - 2 Remove \tilde{t}_j from q by division $q^{\setminus j} \propto q / \tilde{t}_j$.
 - 3 Compute Z_j and \hat{p}_j and find q_{new} by minimizing $\text{KL}(\hat{p}_j || q_{\text{new}})$.
 - 4 Compute and store the new factor $\tilde{t}_j = Z_j q_{\text{new}} / q^{\setminus j}$.
- 3 Evaluate the approximation to the model evidence:

$$p(\mathbf{y}) \approx Z = \int \prod_j \tilde{t}_j(\mathbf{f}) d\mathbf{f}.$$

The FITC prior results in a total cost of $\mathcal{O}(NM^2)$!

Graphical Illustration

Approximates $p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \prod_{j=1}^N t_j(\mathbf{f})$ with $q(\mathbf{f}) \propto t_0(\mathbf{f}) \prod_{j=1}^N \tilde{t}_j(\mathbf{t})$

Graphical Illustration

Approximates $p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \prod_{j=1}^N t_j(\mathbf{f})$ with $q(\mathbf{f}) \propto t_0(\mathbf{f}) \prod_{j=1}^N \tilde{t}_j(\mathbf{t})$

$$p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \quad t_1(\mathbf{f}) \quad t_2(\mathbf{f}) \quad t_3(\mathbf{f}) \quad \approx \quad q(\mathbf{f}) \propto t_0(\mathbf{f}) \quad \tilde{t}_1(\mathbf{f}) \quad \tilde{t}_2(\mathbf{f}) \quad \tilde{t}_3(\mathbf{f})$$


Graphical Illustration

Approximates $p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \prod_{j=1}^N t_j(\mathbf{f})$ with $q(\mathbf{f}) \propto t_0(\mathbf{f}) \prod_{j=1}^N \tilde{t}_j(\mathbf{f})$

$$p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \quad t_1(\mathbf{f}) \quad t_2(\mathbf{f}) \quad t_3(\mathbf{f}) \quad \approx \quad q(\mathbf{f}) \propto t_0(\mathbf{f}) \quad \tilde{t}_1(\mathbf{f}) \quad \tilde{t}_2(\mathbf{f}) \quad \tilde{t}_3(\mathbf{f})$$


The \tilde{t}_j are tuned by minimizing the KL-divergence

$$\text{KL}[\hat{p}_j || q] \quad \text{for } j = 1, \dots, N, \quad \text{where} \quad \hat{p}_j(\mathbf{f}) \propto t_j(\mathbf{f}) \prod_{i \neq j} \tilde{t}_i(\mathbf{f})$$
$$q(\mathbf{f}) \propto \tilde{t}_j(\mathbf{f}) \prod_{i \neq j} \tilde{t}_i(\mathbf{f})$$

Graphical Illustration

Approximates $p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \prod_{j=1}^N t_j(\mathbf{f})$ with $q(\mathbf{f}) \propto t_0(\mathbf{f}) \prod_{j=1}^N \tilde{t}_j(\mathbf{f})$

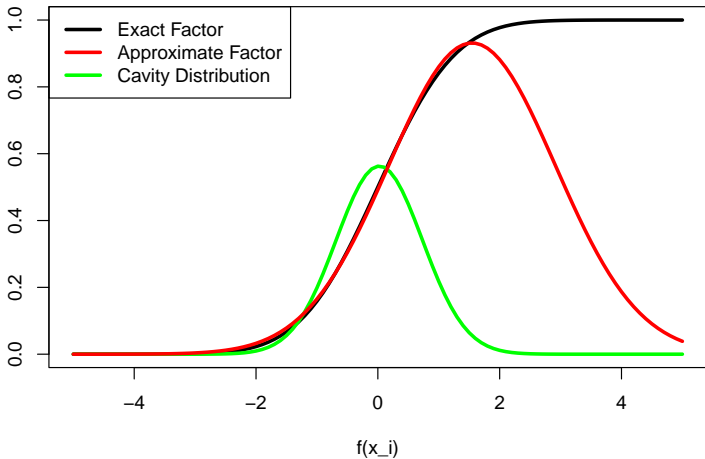
$$p(\mathbf{f}|\mathbf{y}) \propto t_0(\mathbf{f}) \quad t_1(\mathbf{f}) \quad t_2(\mathbf{f}) \quad t_3(\mathbf{f}) \quad \approx \quad q(\mathbf{f}) \propto t_0(\mathbf{f}) \quad \tilde{t}_1(\mathbf{f}) \quad \tilde{t}_2(\mathbf{f}) \quad \tilde{t}_3(\mathbf{f})$$

The \tilde{t}_j are tuned by minimizing the KL-divergence

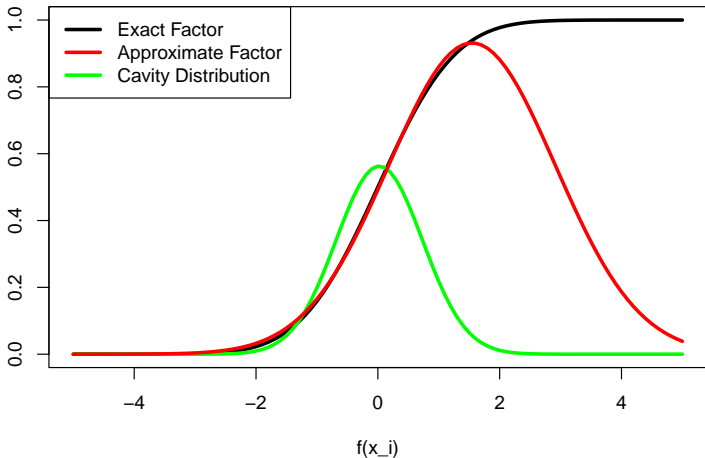
$$\text{KL}[\hat{p}_j || q] \quad \text{for } j = 1, \dots, N, \quad \text{where} \quad \hat{p}_j(\mathbf{f}) \propto t_j(\mathbf{f}) \prod_{i \neq j} \tilde{t}_i(\mathbf{f}) \\ q(\mathbf{f}) \propto \tilde{t}_j(\mathbf{f}) \prod_{i \neq j} \tilde{t}_i(\mathbf{f}) .$$

If the exact factor already belongs to the exponential family it needs not be approximated!

GFITC: Factor Approximation



GFITC: Factor Approximation



The approximate factor is accurate in regions of high posterior probability as indicated by the cavity distribution!

Hyper-parameters and Inducing Points

They are optimized by maximizing the EP estimate of the log-marginal likelihood $\log Z \approx \log p(\mathbf{y})$.

Hyper-parameters and Inducing Points

They are optimized by maximizing the EP estimate of the log-marginal likelihood $\log Z \approx \log p(\mathbf{y})$.

Problem: The parameters θ_i of the approximate factors also depend on the hyper-parameters (including the inducing points)!

Hyper-parameters and Inducing Points

They are optimized by maximizing the EP estimate of the log-marginal likelihood $\log Z \approx \log p(\mathbf{y})$.

Problem: The parameters θ_i of the approximate factors also depend on the hyper-parameters (including the inducing points)!

- Direct dependence of $\log Z$ on the hyper-parameters.

Hyper-parameters and Inducing Points

They are optimized by maximizing the EP estimate of the **log-marginal likelihood** $\log Z \approx \log p(\mathbf{y})$.

Problem: The parameters θ_i of the approximate factors also depend on the hyper-parameters (including the inducing points)!

- Direct dependence of $\log Z$ on the hyper-parameters.
- Indirect dependence of $\log Z$ on the hyper-parameters via each θ_i .

Hyper-parameters and Inducing Points

They are optimized by maximizing the EP estimate of the log-marginal likelihood $\log Z \approx \log p(\mathbf{y})$.

Problem: The parameters θ_i of the approximate factors also depend on the hyper-parameters (including the inducing points)!

- Direct dependence of $\log Z$ on the hyper-parameters.
- Indirect dependence of $\log Z$ on the hyper-parameters via each θ_i .

If EP converges the gradient of $\log Z$ w.r.t. each θ_i is zero, which allows to easily compute the gradients of $\log Z$!

GFITC: Predictions

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

GFITC: Predictions

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

$$p(\mathbf{f}^*, \mathbf{f}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}^* \\ \mathbf{f} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{f}^*\mathbf{f}^*} & \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \\ \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} & \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} \end{bmatrix} \right)$$

GFITC: Predictions

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

$$p(\mathbf{f}^*, \mathbf{f}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}^* \\ \mathbf{f} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{f}^*\mathbf{f}^*} & \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \\ \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} & \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} \end{bmatrix} \right)$$

The conditional $p(\mathbf{f}^*|\mathbf{f})$ is:

$$p(\mathbf{f}^*|\mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} \right)$$

GFITC: Predictions

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

$$p(\mathbf{f}^*, \mathbf{f}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}^* \\ \mathbf{f} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{f}^*\mathbf{f}^*} & \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \\ \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} & \tilde{\Sigma}_{\mathbf{f}\mathbf{f}} \end{bmatrix} \right)$$

The conditional $p(\mathbf{f}^*|\mathbf{f})$ is:

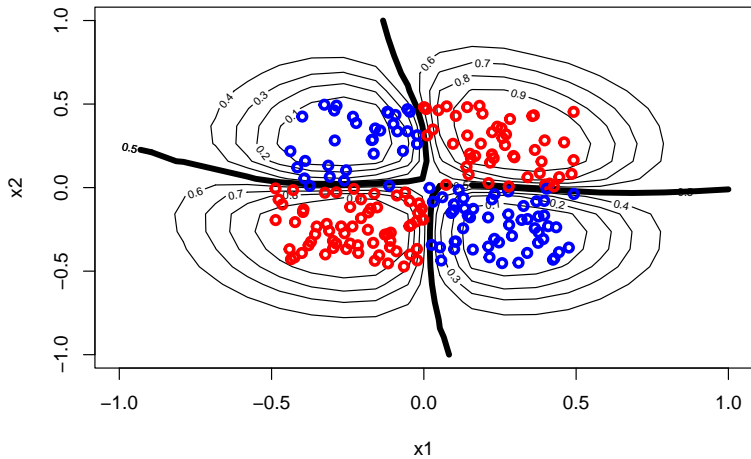
$$p(\mathbf{f}^*|\mathbf{f}) = \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{f}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \tilde{\Sigma}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} \right)$$

After marginalizing \mathbf{f} w.r.t. $q(\mathbf{f})$, we obtain the predictive distribution:

$$\begin{aligned} p(\mathbf{f}^*|\mathbf{y}) &= \int p(\mathbf{f}^*|\mathbf{f})q(\mathbf{f})d\mathbf{f} \\ &= \mathcal{N} \left(\mathbf{f}^* \mid \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \left(\tilde{\Sigma}_{\mathbf{f}\mathbf{f}} + \tilde{\mathbf{\Pi}} \right)^{-1} \tilde{\mathbf{y}}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \mathbf{Q}_{\mathbf{f}^*\mathbf{f}} \left(\tilde{\Sigma}_{\mathbf{f}\mathbf{f}} + \tilde{\mathbf{\Pi}} \right)^{-1} \mathbf{Q}_{\mathbf{f}\mathbf{f}^*} \right) \end{aligned}$$

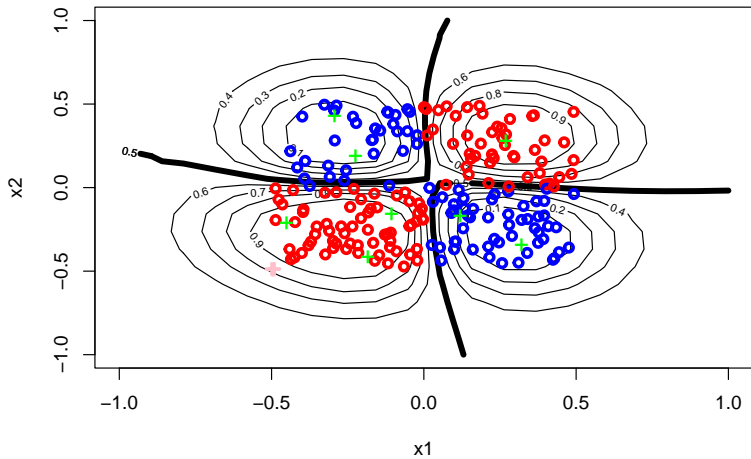
GFITC: Illustrative Example

Full GP + EP



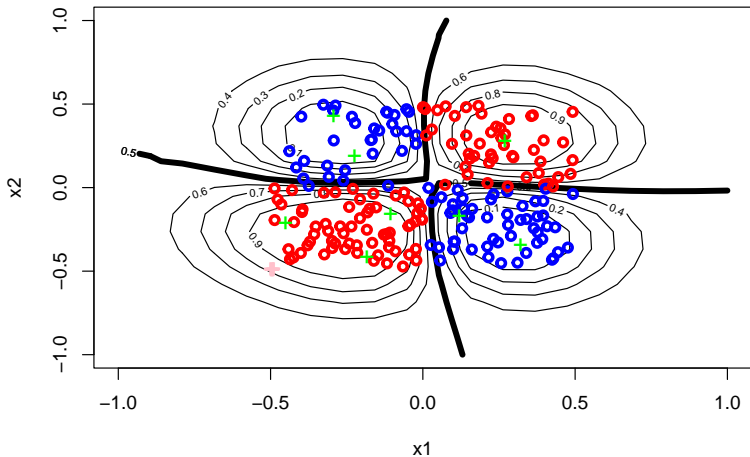
GFITC: Illustrative Example

GFITC (M=10)



GFITC: Illustrative Example

GFITC (M=10)



The inducing points spread across the input space!

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy (VFE) Method:

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy (VFE) Method:

- Keeps the GP prior intact and does not introduce any simplification!

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy (VFE) Method:

- Keeps the GP prior intact and does not introduce any simplification!
- Carries out approximate inference to approximate the GP posterior.

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy (VFE) Method:

- Keeps the GP prior intact and does not introduce any simplification!
- Carries out approximate inference to approximate the GP posterior.
- The particular approximate distribution q results in cost $\mathcal{O}(NM^2)$.

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy (VFE) Method:

- Keeps the GP prior intact and does not introduce any simplification!
- Carries out approximate inference to approximate the GP posterior.
- The particular approximate distribution q results in cost $\mathcal{O}(NM^2)$.
- Variational inference is used to tune q .

Variational Free Energy

Previous methods approximate the GP prior using a low rank approximation of Σ , resulting in a cost $\mathcal{O}(NM^2)$.

Variational Free Energy (VFE) Method:

- Keeps the GP prior intact and does not introduce any simplification!
- Carries out approximate inference to approximate the GP posterior.
- The particular approximate distribution q results in cost $\mathcal{O}(NM^2)$.
- Variational inference is used to tune q .

Since the GP prior is not changed it tends to perform better than the previous methods!

Variational Inference

Adjust the parameters of q to match p by minimizing $\text{KL}(q|p) \geq 0$.

Variational Inference

Adjust the parameters of q to match p by minimizing $\text{KL}(q|p) \geq 0$.

$$\text{KL}(q|p) = 0 \iff q(\mathbf{f}) = p(\mathbf{f})$$

Variational Inference

Adjust the parameters of q to match p by minimizing $\text{KL}(q|p) \geq 0$.

$$\text{KL}(q|p) = 0 \iff q(\mathbf{f}) = p(\mathbf{f})$$

The expression for the KL divergence between q and p is:

$$\int q(\mathbf{f}) \log \frac{q(\mathbf{f})}{p(\mathbf{f})} d\mathbf{f} \geq 0$$

Variational Inference

Adjust the parameters of q to match p by minimizing $\text{KL}(q|p) \geq 0$.

$$\text{KL}(q|p) = 0 \iff q(\mathbf{f}) = p(\mathbf{f})$$

The expression for the KL divergence between q and p is:

$$\int q(\mathbf{f}) \log \frac{q(\mathbf{f})}{p(\mathbf{f})} d\mathbf{f} \geq 0$$

$\text{KL}(q|p)$ depends on p , which is assumed to be intractable!

Variational Inference

Adjust the parameters of q to match p by minimizing $\text{KL}(q|p) \geq 0$.

$$\text{KL}(q|p) = 0 \iff q(\mathbf{f}) = p(\mathbf{f})$$

The expression for the KL divergence between q and p is:

$$\int q(\mathbf{f}) \log \frac{q(\mathbf{f})}{p(\mathbf{f})} d\mathbf{f} \geq 0$$

KL($q|p$) depends on p , which is assumed to be intractable!

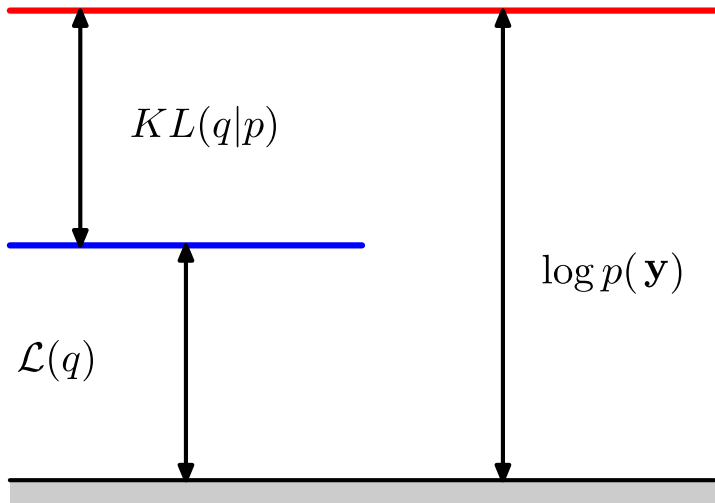
Let the target be $p(\mathbf{f}|\mathbf{y})$. Consider the decomposition of $p(\mathbf{y})$:

$$\log p(\mathbf{y}) = \mathcal{L}(q) + \text{KL}(q|p),$$

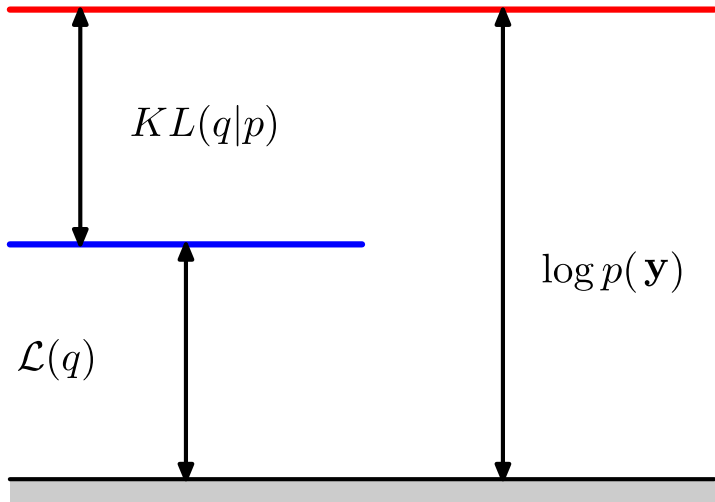
where

$$\mathcal{L}(q) = \int q(\mathbf{f}) \log \frac{p(\mathbf{f}, \mathbf{y})}{q(\mathbf{f})} d\mathbf{f}, \quad \text{KL}(q|p) = \int q(\mathbf{f}) \log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} d\mathbf{f}.$$

Decomposition of the Marginal Likelihood



Decomposition of the Marginal Likelihood

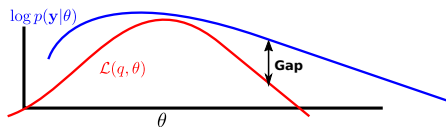


$\mathcal{L}(q)$ can be used to approximate $\log p(\mathbf{y})$ if $KL(q|p)$ is small!

Variational Free Energy (VFE)

Lower bound the log-likelihood:

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) d\mathbf{f} d\mathbf{u}$$

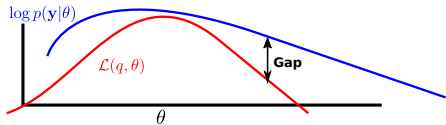


Variational Free Energy (VFE)

Lower bound the log-likelihood:

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) d\mathbf{f} d\mathbf{u}$$

$$= \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u}$$

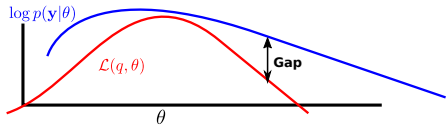


Variational Free Energy (VFE)

Lower bound the log-likelihood:

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) d\mathbf{f} d\mathbf{u}$$

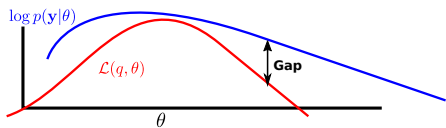
$$= \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \geq \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \equiv \mathcal{L}(q, \theta)$$



Variational Free Energy (VFE)

Lower bound the log-likelihood:

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) d\mathbf{f} d\mathbf{u}$$



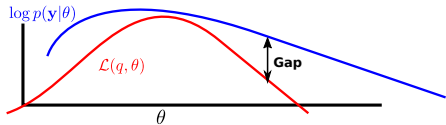
$$= \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \geq \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \equiv \mathcal{L}(q, \theta)$$

$$\mathcal{L}(q, \theta) = \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} = \log p(\mathbf{y}|\theta) - \mathbf{KL}[q(\mathbf{f}, \mathbf{u})|p(\mathbf{f}, \mathbf{u}|\mathbf{y})]$$

Variational Free Energy (VFE)

Lower bound the log-likelihood:

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) d\mathbf{f} d\mathbf{u}$$



$$= \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \geq \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \equiv \mathcal{L}(q, \theta)$$

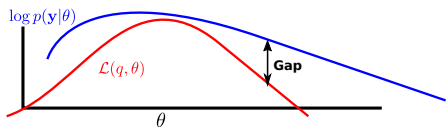
$$\mathcal{L}(q, \theta) = \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} = \log p(\mathbf{y}|\theta) - \mathbf{KL}[q(\mathbf{f}, \mathbf{u})|p(\mathbf{f}, \mathbf{u}|\mathbf{y})]$$

$\mathbf{KL} \equiv$ Kullback-Leibler divergence

Variational Free Energy (VFE)

Lower bound the log-likelihood:

$$\log p(\mathbf{y}|\theta) = \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) d\mathbf{f} d\mathbf{u}$$



$$= \log \int p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \geq \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \equiv \mathcal{L}(q, \theta)$$

$$\mathcal{L}(q, \theta) = \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} = \log p(\mathbf{y}|\theta) - \mathbf{KL}[q(\mathbf{f}, \mathbf{u})|p(\mathbf{f}, \mathbf{u}|\mathbf{y})]$$

$\mathbf{KL} \equiv$ Kullback-Leibler divergence

By maximizing $\mathcal{L}(q, \theta)$ w.r.t. q we are enforcing that $q(\mathbf{f}, \mathbf{u})$ looks similar to $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$ in terms of the KL!

Variational Free Energy (VFE)

Consider the following approximate distribution:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u}) q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u}) \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$$

Variational Free Energy (VFE)

Consider the following approximate distribution:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u}) q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u}) \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$$

- Fixed
- Tunable

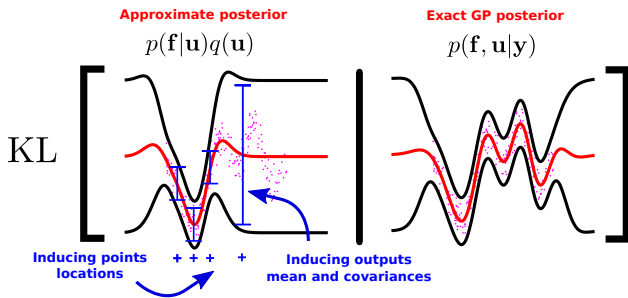


Variational Free Energy (VFE)

Consider the following approximate distribution:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u}) q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u}) \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$$

- Fixed
- Tunable

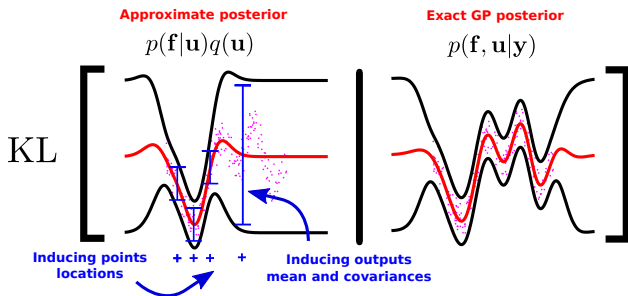


Variational Free Energy (VFE)

Consider the following approximate distribution:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u}) q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u}) \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$$

- Fixed
- Tunable



The inducing points are now parameters of the approx. dist. $q!$

Variational Free Energy (VFE)

Plugging $q(\mathbf{f}, \mathbf{u})$ into the lower bound we have:

Variational Free Energy (VFE)

Plugging $q(\mathbf{f}, \mathbf{u})$ into the lower bound we have:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u} | \theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \\ &= \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \frac{p(\mathbf{y} | \mathbf{f}, \theta) p(\mathbf{f} | \mathbf{u}) p(\mathbf{u})}{p(\mathbf{f} | \mathbf{u}) q(\mathbf{u})} d\mathbf{f} d\mathbf{u}\end{aligned}$$

Variational Free Energy (VFE)

Plugging $q(\mathbf{f}, \mathbf{u})$ into the lower bound we have:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u} | \theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \\ &= \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \frac{p(\mathbf{y} | \mathbf{f}, \theta) \cancel{p(\mathbf{f} | \mathbf{u})} p(\mathbf{u})}{\cancel{p(\mathbf{f} | \mathbf{u})} q(\mathbf{u})} d\mathbf{f} d\mathbf{u}\end{aligned}$$

Variational Free Energy (VFE)

Plugging $q(\mathbf{f}, \mathbf{u})$ into the lower bound we have:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u} | \theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \\ &= \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \frac{p(\mathbf{y} | \mathbf{f}, \theta) \cancel{p(\mathbf{f} | \mathbf{u})} p(\mathbf{u})}{\cancel{p(\mathbf{f} | \mathbf{u})} q(\mathbf{u})} d\mathbf{f} d\mathbf{u}\end{aligned}$$

$$\mathcal{L}(q, \theta) = \mathbb{E}_{q(\mathbf{f})}[\log p(\mathbf{y} | \mathbf{f}, \theta)] - \text{KL}[q(\mathbf{u}) | p(\mathbf{u})]$$

- Mean squared prediction error
- KL between Gaussians

Variational Free Energy (VFE)

Plugging $q(\mathbf{f}, \mathbf{u})$ into the lower bound we have:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u} | \theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \\ &= \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \frac{p(\mathbf{y} | \mathbf{f}, \theta) \cancel{p(\mathbf{f} | \mathbf{u})} p(\mathbf{u})}{\cancel{p(\mathbf{f} | \mathbf{u})} q(\mathbf{u})} d\mathbf{f} d\mathbf{u}\end{aligned}$$

$$\mathcal{L}(q, \theta) = \mathbb{E}_{q(\mathbf{f})}[\log p(\mathbf{y} | \mathbf{f}, \theta)] - \text{KL}[q(\mathbf{u}) | p(\mathbf{u})]$$

- Mean squared prediction error
- KL between Gaussians
- No change in the model is made and the cost is in $\mathcal{O}(M^2N)$!

Variational Free Energy (VFE)

Plugging $q(\mathbf{f}, \mathbf{u})$ into the lower bound we have:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \int q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u} | \theta)}{q(\mathbf{f}, \mathbf{u})} d\mathbf{f} d\mathbf{u} \\ &= \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \frac{p(\mathbf{y} | \mathbf{f}, \theta) p(\mathbf{f} | \mathbf{u}) p(\mathbf{u})}{p(\mathbf{f} | \mathbf{u}) q(\mathbf{u})} d\mathbf{f} d\mathbf{u}\end{aligned}$$

$$\mathcal{L}(q, \theta) = \mathbb{E}_{q(\mathbf{f})}[\log p(\mathbf{y} | \mathbf{f}, \theta)] - \text{KL}[q(\mathbf{u}) | p(\mathbf{u})]$$

- Mean squared prediction error
- KL between Gaussians
- No change in the model is made and the cost is in $\mathcal{O}(M^2N)$!
- Predictions are made using $p(\mathbf{f}^* | \mathbf{u}) q(\mathbf{u})$ marginalizing out \mathbf{u} .

VFE: Predictions for Test Instances

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

VFE: Predictions for Test Instances

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

$$p(\mathbf{f}^*, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}^* \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{f}^*\mathbf{f}^*} & \Sigma_{\mathbf{f}^*\mathbf{u}} \\ \Sigma_{\mathbf{u}\mathbf{f}^*} & \Sigma_{\mathbf{u}\mathbf{u}} \end{bmatrix} \right)$$

VFE: Predictions for Test Instances

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

$$p(\mathbf{f}^*, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}^* \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{f}^*\mathbf{f}^*} & \Sigma_{\mathbf{f}^*\mathbf{u}} \\ \Sigma_{\mathbf{u}\mathbf{f}^*} & \Sigma_{\mathbf{u}\mathbf{u}} \end{bmatrix} \right)$$

The conditional $p(\mathbf{f}^*|\mathbf{u})$ is:

$$p(\mathbf{f}^*|\mathbf{u}) = \mathcal{N}(\mathbf{f}^* | \Sigma_{\mathbf{f}^*\mathbf{u}} \Sigma_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{u}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \Sigma_{\mathbf{f}^*\mathbf{u}} \Sigma_{\mathbf{u}\mathbf{u}}^{-1} \Sigma_{\mathbf{u}\mathbf{f}^*})$$

VFE: Predictions for Test Instances

We want to compute the value of \mathbf{f}^* at a new \mathbf{x}^* :

$$p(\mathbf{f}^*, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}^* \\ \mathbf{u} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{f}^*\mathbf{f}^*} & \Sigma_{\mathbf{f}^*\mathbf{u}} \\ \Sigma_{\mathbf{u}\mathbf{f}^*} & \Sigma_{\mathbf{u}\mathbf{u}} \end{bmatrix} \right)$$

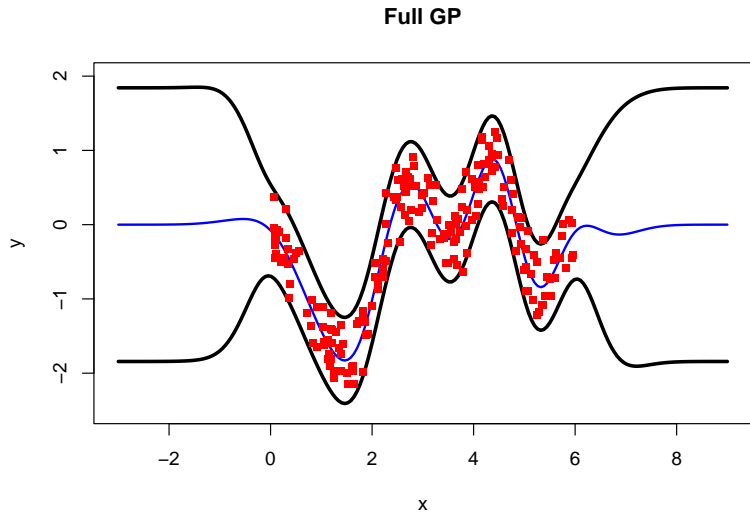
The conditional $p(\mathbf{f}^*|\mathbf{u})$ is:

$$p(\mathbf{f}^*|\mathbf{u}) = \mathcal{N}(\mathbf{f}^* | \Sigma_{\mathbf{f}^*\mathbf{u}}\Sigma_{\mathbf{u}\mathbf{u}}^{-1}\mathbf{u}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \Sigma_{\mathbf{f}^*\mathbf{u}}\Sigma_{\mathbf{u}\mathbf{u}}^{-1}\Sigma_{\mathbf{u}\mathbf{f}^*})$$

After marginalizing \mathbf{u} w.r.t. $q(\mathbf{u})$, we obtain the predictive distribution:

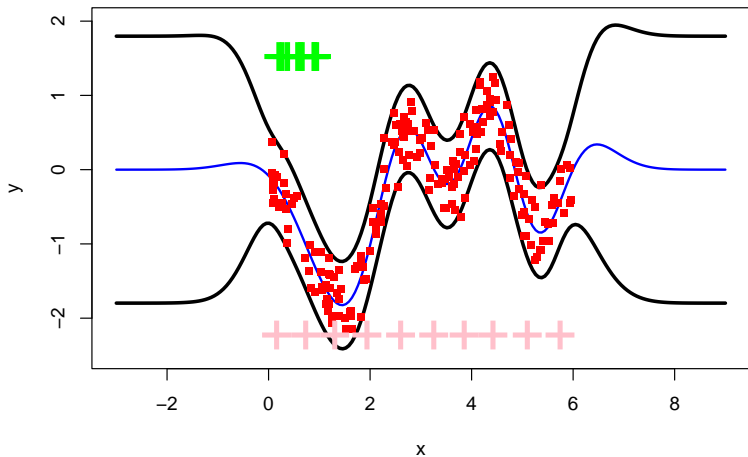
$$\begin{aligned} p(\mathbf{f}^*|\mathbf{y}) &= \int p(\mathbf{f}^*|\mathbf{u})q(\mathbf{u})d\mathbf{u} \\ &= \mathcal{N}(\mathbf{f}^* | \Sigma_{\mathbf{f}^*\mathbf{u}}\Sigma_{\mathbf{u}\mathbf{u}}^{-1}\mathbf{m}, \Sigma_{\mathbf{f}^*\mathbf{f}^*} - \Sigma_{\mathbf{f}^*\mathbf{u}}(\Sigma_{\mathbf{u}\mathbf{u}}^{-1} - \Sigma_{\mathbf{u}\mathbf{u}}^{-1}\mathbf{S}\Sigma_{\mathbf{u}\mathbf{u}}^{-1})\Sigma_{\mathbf{u}\mathbf{f}^*}) \end{aligned}$$

VFE: Illustrative Example



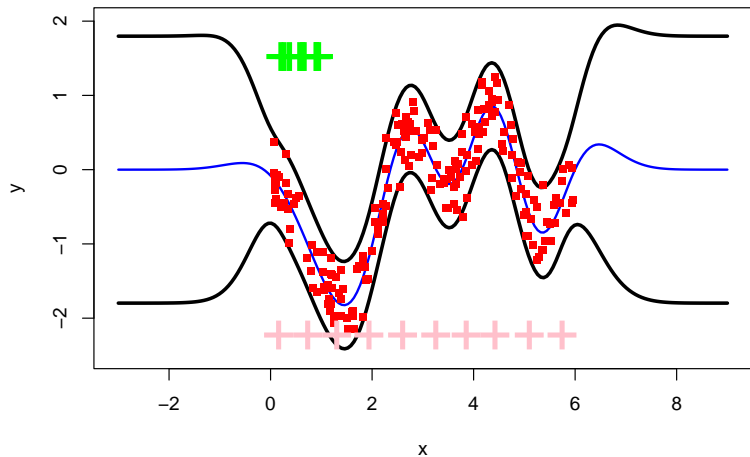
VFE: Illustrative Example

VFE ($M = 10$)



VFE: Illustrative Example

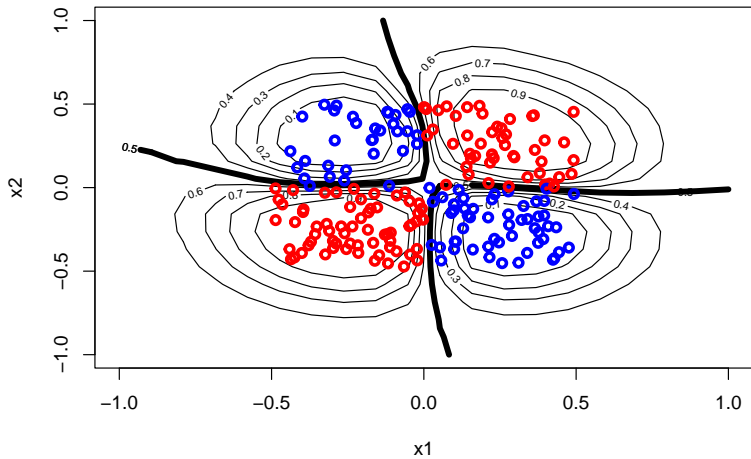
VFE ($M = 10$)



The inducing points cover the regions where the function changes!

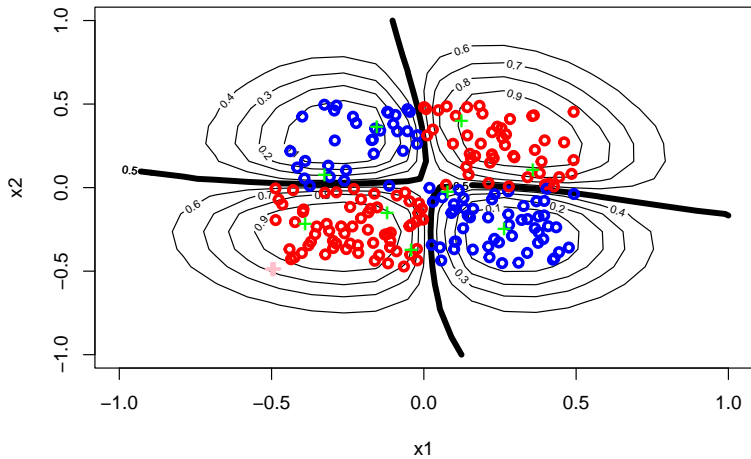
VFE: Illustrative Classification Example

Full GP + EP

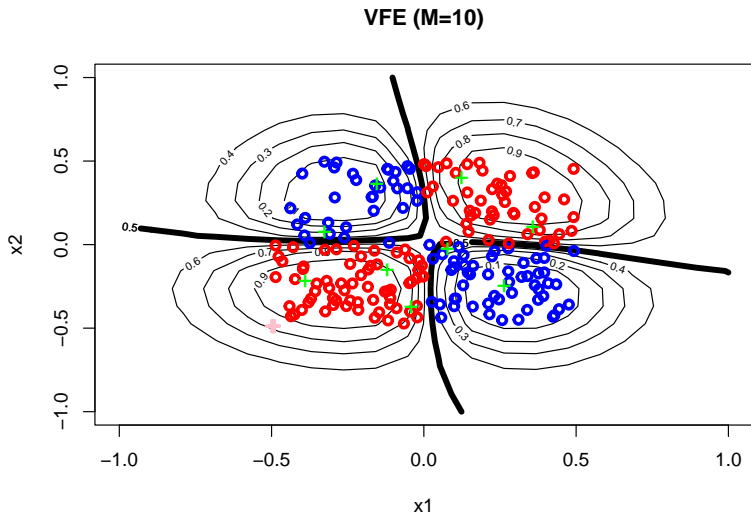


VFE: Illustrative Classification Example

VFE (M=10)



VFE: Illustrative Classification Example



The inducing points spread across the input space!

FITC vs. VFE

Two approaches:

FITC vs. VFE

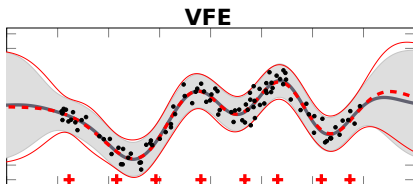
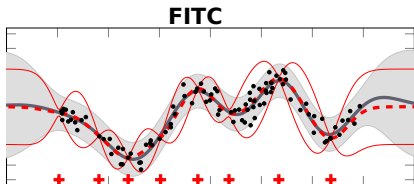
Two approaches:

- FITC: optimize the marginal likelihood of an approximate GP model.
- VFE: maximize fidelity to the original exact GP.

FITC vs. VFE

Two approaches:

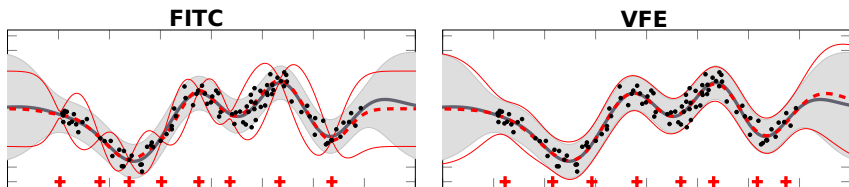
- FITC: optimize the marginal likelihood of an approximate GP model.
- VFE: maximize fidelity to the original exact GP.



FITC vs. VFE

Two approaches:

- FITC: optimize the marginal likelihood of an approximate GP model.
- VFE: maximize fidelity to the original exact GP.



- FITC: less local optima and easier to optimize, also less accurate.
- VFE: more accurate, more local optima, more difficult to optimize.

(Bui et al., 2017) (Bauer et al., 2016)

Run the cells of the notebook to fit a sparse GP using the VFE approximation and complete task 6!

Whitened Parameterization for VFE

Alternative VFE objective expected to be easier to optimize!

Whitened Parameterization for VFE

Alternative VFE objective expected to be easier to optimize!

Instead of making inference about \mathbf{u} , the whitened VFE objective makes inference about:

$$\mathbf{e} \quad \text{such that} \quad \mathbf{u} = \mathbf{L}\mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

with \mathbf{u} the latent process values at the inducing points and $\mathbf{L}^T\mathbf{L} = \Sigma_{\mathbf{u}\mathbf{u}}$.

Whitened Parameterization for VFE

Alternative VFE objective expected to be easier to optimize!

Instead of making inference about \mathbf{u} , the whitened VFE objective makes inference about:

$$\mathbf{e} \quad \text{such that} \quad \mathbf{u} = \mathbf{L}\mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

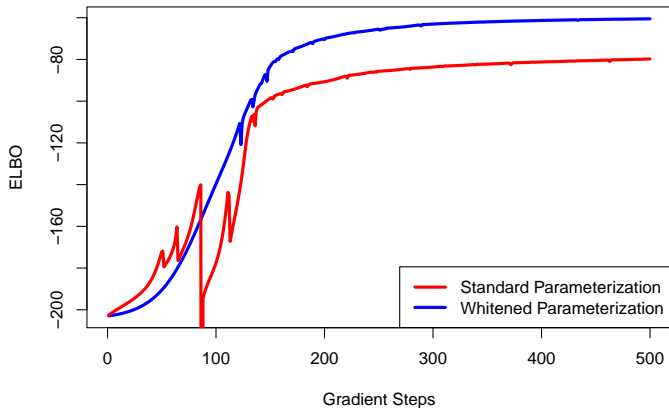
with \mathbf{u} the latent process values at the inducing points and $\mathbf{L}^T \mathbf{L} = \Sigma_{\mathbf{u}\mathbf{u}}$.

The VFE objective becomes:

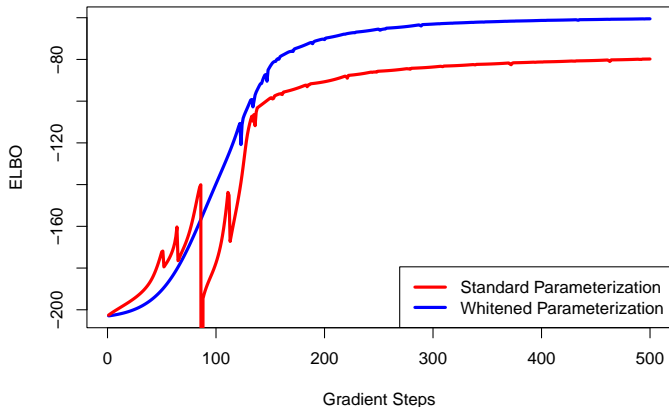
$$\sum_{i=1}^N \mathbb{E}_{q(\mathbf{e})p(f(\mathbf{x}_i)|\mathbf{e})} [\log p(y_i|f(\mathbf{x}_i))] - \text{KL}(q(\mathbf{e})|\mathcal{N}(\mathbf{0}, \mathbf{I})),$$

with $p(f(\mathbf{x}_i)|\mathbf{e})$ using the covariances between $f(\mathbf{x}_i)$ and \mathbf{e} .

Whitened Parameterization: Illustrative Example



Whitened Parameterization: Illustrative Example



Whitening significantly improves convergence!

**Run the cells of the notebook to fit a sparse GP
using the VFE approximation with whitening!**

Natural Gradient Ascent

Gradient ascent moves in the direction of the gradient $\nabla_{\xi} \mathcal{L}(\xi)$.

Natural Gradient Ascent

Gradient ascent moves in the direction of the gradient $\nabla_{\xi} \mathcal{L}(\xi)$.

Formally:

$$\nabla_{\xi} \mathcal{L}(\xi) \propto \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max_{\mathbf{d} \text{ s.t. } \|\mathbf{d}\| \leq \epsilon} \mathcal{L}(\xi + \mathbf{d})$$

Natural Gradient Ascent

Gradient ascent moves in the direction of the gradient $\nabla_{\xi} \mathcal{L}(\xi)$.

Formally:

$$\nabla_{\xi} \mathcal{L}(\xi) \propto \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max_{\mathbf{d} \text{ s.t. } \|\mathbf{d}\| \leq \epsilon} \mathcal{L}(\xi + \mathbf{d})$$

The steepest ascent direction picks \mathbf{d} in the ϵ -vicinity of ξ (measured by the Euclidean norm) that maximizes $\mathcal{L}(\cdot)$.

Natural Gradient Ascent

Gradient ascent moves in the direction of the gradient $\nabla_{\xi} \mathcal{L}(\xi)$.

Formally:

$$\nabla_{\xi} \mathcal{L}(\xi) \propto \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max_{\mathbf{d} \text{ s.t. } \|\mathbf{d}\| \leq \epsilon} \mathcal{L}(\xi + \mathbf{d})$$

The steepest ascent direction picks \mathbf{d} in the ϵ -vicinity of ξ (measured by the Euclidean norm) that maximizes $\mathcal{L}(\cdot)$.

If ξ represents the parameters of probability distributions, the Euclidean norm may be problematic!

Illustration with Two Gaussians

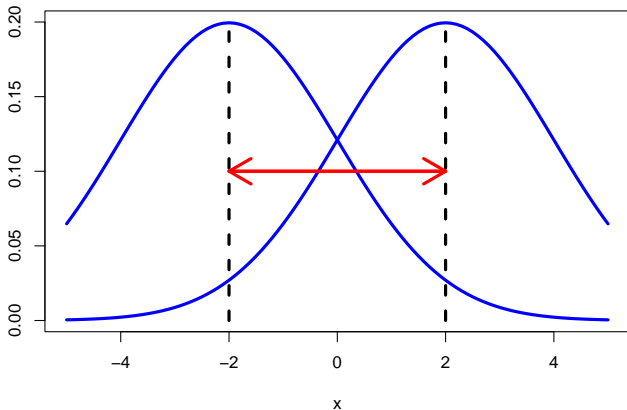


Illustration with Two Gaussians

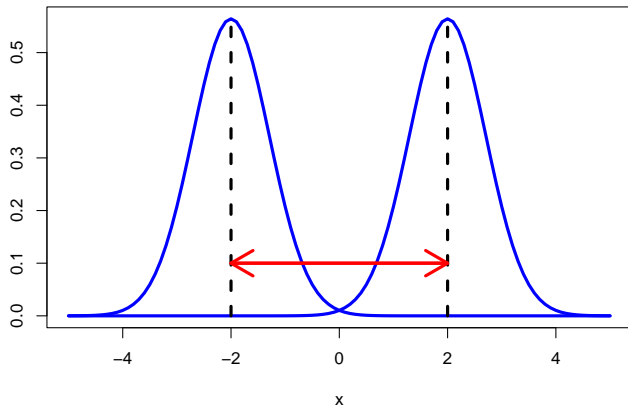
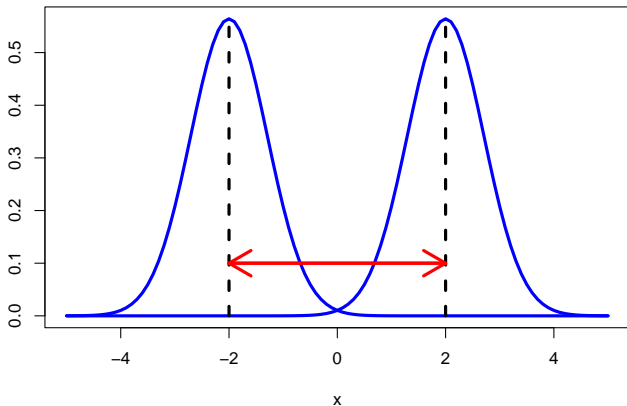
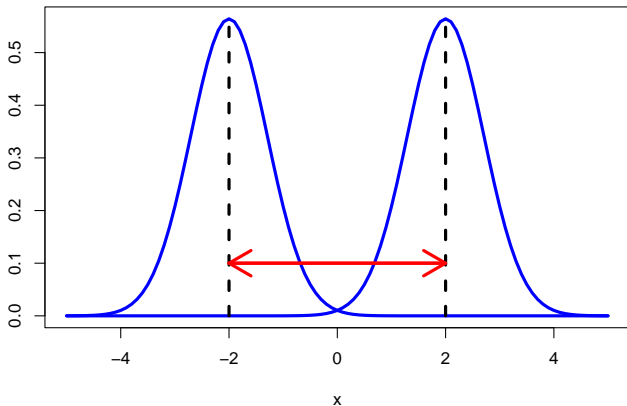


Illustration with Two Gaussians



The Euclidean distance between parameters is 4 in both cases!

Illustration with Two Gaussians



A better alternative is the KL-divergence between distributions!

Natural Gradient Ascent (efficient computation)

Considers the KL-divergence as a norm:

$$\nabla_{\xi} \mathcal{L}(\xi) \mathbf{F}_{\xi}^{-1} \propto \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max_{\mathbf{d} \text{ s.t. } \text{KL}[q(\mathbf{u}|\xi)|q(\mathbf{u}|\xi+\mathbf{d})] \leq \epsilon} \mathcal{L}(\xi + \mathbf{d})$$

with \mathbf{F}_{ξ} the Fisher information of q :

$$\mathbf{F}_{\xi} = -\mathbb{E}_{q(\mathbf{u}|\xi)} [\nabla_{\xi}^2 \log q(\mathbf{u}|\xi)]$$

Natural Gradient Ascent (efficient computation)

Considers the KL-divergence as a norm:

$$\nabla_{\xi} \mathcal{L}(\xi) \mathbf{F}_{\xi}^{-1} \propto \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max_{\mathbf{d} \text{ s.t. } \text{KL}[q(\mathbf{u}|\xi)|q(\mathbf{u}|\xi+\mathbf{d})] \leq \epsilon} \mathcal{L}(\xi + \mathbf{d})$$

with \mathbf{F}_{ξ} the Fisher information of q :

$$\mathbf{F}_{\xi} = -\mathbb{E}_{q(\mathbf{u}|\xi)} [\nabla_{\xi}^2 \log q(\mathbf{u}|\xi)]$$

Let η and θ be the natural and expectation parameters of q , respectively:

$$\mathbf{F}_{\eta} = \frac{\partial \theta}{\partial \eta}, \quad \mathbf{F}_{\xi} = \left(\frac{\partial \eta}{\partial \xi} \right)^{\top} \frac{\partial \theta}{\partial \eta} \frac{\partial \eta}{\partial \xi} = \left(\frac{\partial \eta}{\partial \xi} \right)^{\top} \frac{\partial \theta}{\partial \xi}.$$

Natural Gradient Ascent (efficient computation)

Considers the KL-divergence as a norm:

$$\nabla_{\xi} \mathcal{L}(\xi) \mathbf{F}_{\xi}^{-1} \propto \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max_{\mathbf{d} \text{ s.t. } \text{KL}[q(\mathbf{u}|\xi)|q(\mathbf{u}|\xi+\mathbf{d})] \leq \epsilon} \mathcal{L}(\xi + \mathbf{d})$$

with \mathbf{F}_{ξ} the Fisher information of q :

$$\mathbf{F}_{\xi} = -\mathbb{E}_{q(\mathbf{u}|\xi)} [\nabla_{\xi}^2 \log q(\mathbf{u}|\xi)]$$

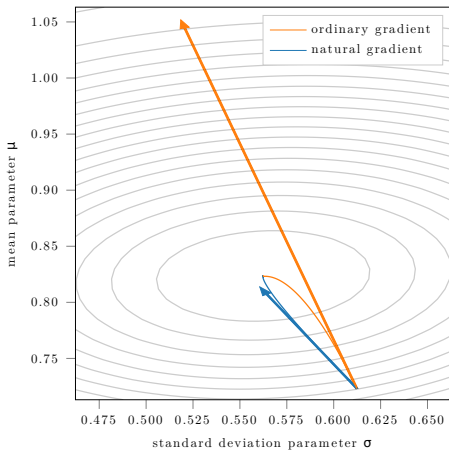
Let η and θ be the natural and expectation parameters of q , respectively:

$$\mathbf{F}_{\eta} = \frac{\partial \theta}{\partial \eta}, \quad \mathbf{F}_{\xi} = \left(\frac{\partial \eta}{\partial \xi} \right)^{\top} \frac{\partial \theta}{\partial \eta} \frac{\partial \eta}{\partial \xi} = \left(\frac{\partial \eta}{\partial \xi} \right)^{\top} \frac{\partial \theta}{\partial \xi}.$$

Thus, we obtain natural gradients in terms of jacobian vector products:

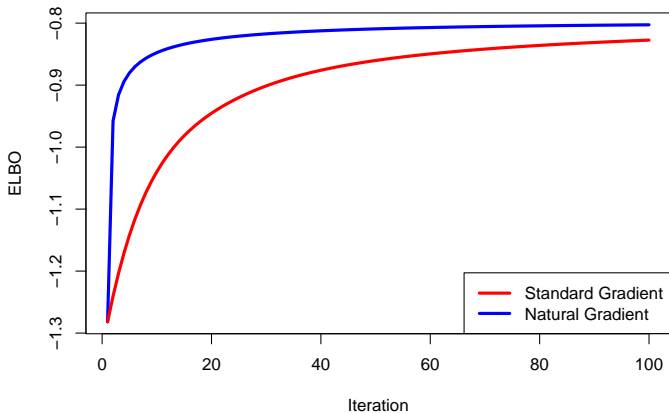
$$\nabla_{\xi} \mathcal{L}(\xi) \mathbf{F}_{\xi}^{-1} = \nabla_{\xi} \mathcal{L}(\xi) \frac{\partial \xi}{\partial \theta} \left(\frac{\partial \xi}{\partial \eta} \right)^{\top} = \frac{\partial \mathcal{L}}{\partial \theta} \left(\frac{\partial \xi}{\partial \eta} \right)^{\top}.$$

Natural Gradient Ascent

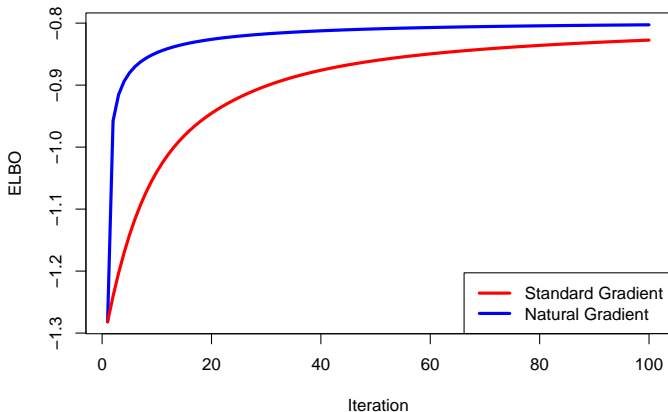


(Salimbeni et al., 2018)

Natural Gradient: Illustrative Example



Natural Gradient: Illustrative Example



The natural gradient achieves a faster convergence!

**Run the cells of the notebook to fit a sparse GP
using the VFE approximation with natural
gradients!**

GPs for Big Data

Can we further improve the computational cost in $\mathcal{O}(NM^2)$?

GPs for Big Data

Can we further improve the computational cost in $\mathcal{O}(NM^2)$?

Minibatch training in NN allows to scale to massive datasets!

GPs for Big Data

Can we further improve the computational cost in $\mathcal{O}(NM^2)$?

Minibatch training in NN allows to scale to massive datasets!

Straight forward to do that in the VFE approach:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \mathbb{E}_{q(\mathbf{f})}[\log p(\mathbf{y}|\mathbf{f}, \theta)] - \mathbf{KL}[q(\mathbf{u})|\rho(\mathbf{u})] \\ &= \sum_{i=1}^N \mathbb{E}_{q(f_i)}[\log p(y_i|f_i, \theta)] - \mathbf{KL}[q(\mathbf{u})|\rho(\mathbf{u})] \\ &\approx \frac{N}{B} \sum_{i \in \mathcal{B}} \mathbb{E}_{q(f_i)}[\log p(y_i|f_i, \theta)] - \mathbf{KL}[q(\mathbf{u})|\rho(\mathbf{u})]\end{aligned}$$

GPs for Big Data

Can we further improve the computational cost in $\mathcal{O}(NM^2)$?

Minibatch training in NN allows to scale to massive datasets!

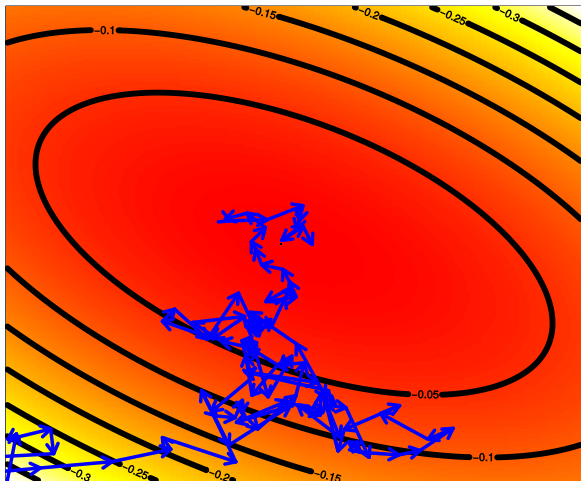
Straight forward to do that in the VFE approach:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \mathbb{E}_{q(\mathbf{f})}[\log p(\mathbf{y}|\mathbf{f}, \theta)] - \mathbf{KL}[q(\mathbf{u})|\rho(\mathbf{u})] \\ &= \sum_{i=1}^N \mathbb{E}_{q(f_i)}[\log p(y_i|f_i, \theta)] - \mathbf{KL}[q(\mathbf{u})|\rho(\mathbf{u})] \\ &\approx \frac{N}{B} \sum_{i \in \mathcal{B}} \mathbb{E}_{q(f_i)}[\log p(y_i|f_i, \theta)] - \mathbf{KL}[q(\mathbf{u})|\rho(\mathbf{u})]\end{aligned}$$

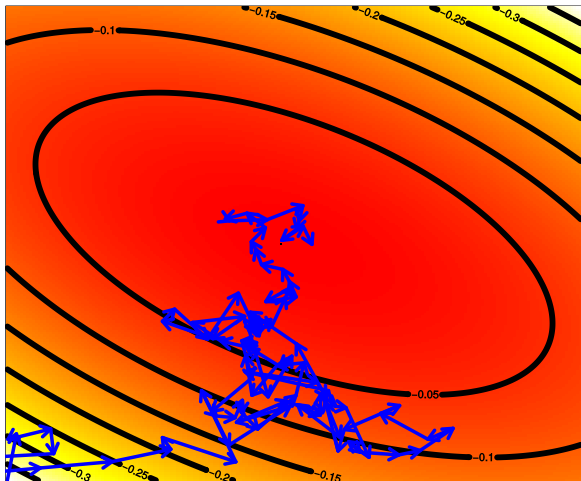
The training cost goes down to $\mathcal{O}(M^3)$ which allows to address datasets with millions of instances!

(Hensman et al., 2013)

GPs for Big Data

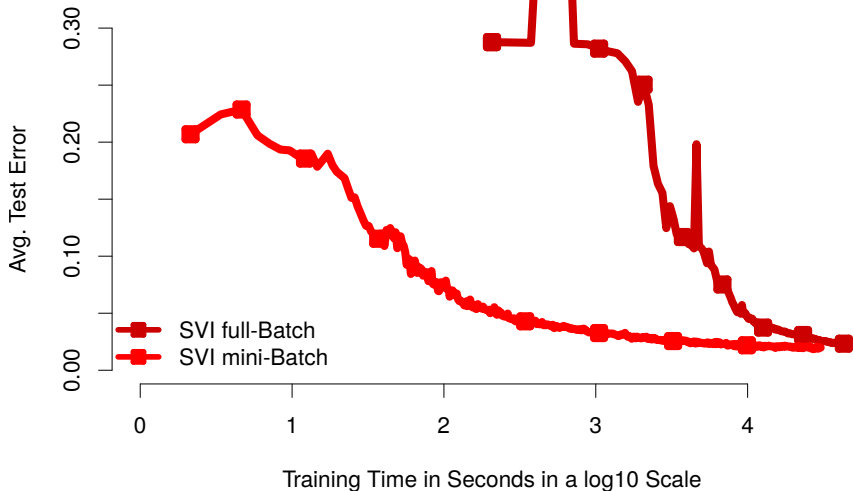


GPs for Big Data



To converge to a local neighborhood of the optimum stochastic methods require an estimate of the gradient which can be very cheap!

GPs for Big Data



(Hernández-Lobato, 2015)

Summary about VFE

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.

Summary about VFE

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.

Summary about VFE

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.
- Guaranteed to be exact if $M = N$ and the inducing points are not optimized and located at the training points.

Summary about VFE

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.
- Guaranteed to be exact if $M = N$ and the inducing points are not optimized and located at the training points.
- It does not change the model. It relies on a particular posterior approximation that speeds-up the computations.

Summary about VFE

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.
- Guaranteed to be exact if $M = N$ and the inducing points are not optimized and located at the training points.
- It does not change the model. It relies on a particular posterior approximation that speeds-up the computations.
- It allows for minibatch training which reduces the cost to $\mathcal{O}(M^3)$.

Summary about VFE

- Reduces the cost to $\mathcal{O}(M^2N)$ with $M \ll N$.
- The optimized inducing points spread over the input space where the latent function changes.
- Guaranteed to be exact if $M = N$ and the inducing points are not optimized and located at the training points.
- It does not change the model. It relies on a particular posterior approximation that speeds-up the computations.
- It allows for minibatch training which reduces the cost to $\mathcal{O}(M^3)$.
- The objective is prone to local optima and difficult to optimize.

Run the cells of the notebook to fit a sparse GP using the VFE approximation with mini-batches!

Run the cells of the notebook to fit a sparse GP using the VFE approximation for classification and carry out task 7!

Sparse GP Conclusions

- Exact GPs have an $\mathcal{O}(N^3)$ computational cost, making them feasible on small datasets with a few thousand instances only.

Sparse GP Conclusions

- Exact GPs have an $\mathcal{O}(N^3)$ computational cost, making them feasible on small datasets with a few thousand instances only.
- Sparse GPs provide an approximate solution with a smaller computational cost that is $\mathcal{O}(NM^2)$ with $M \ll N$.

Sparse GP Conclusions

- Exact GPs have an $\mathcal{O}(N^3)$ computational cost, making them feasible on small datasets with a few thousand instances only.
- Sparse GPs provide an approximate solution with a smaller computational cost that is $\mathcal{O}(NM^2)$ with $M \ll N$.
- The non-parametric property of GP is lost when using sparse approximations. They are no longer more flexible with more data.

Sparse GP Conclusions

- Exact GPs have an $\mathcal{O}(N^3)$ computational cost, making them feasible on small datasets with a few thousand instances only.
- Sparse GPs provide an approximate solution with a smaller computational cost that is $\mathcal{O}(NM^2)$ with $M \ll N$.
- The non-parametric property of GP is lost when using sparse approximations. They are no longer more flexible with more data.
- The methods that approximate the GP prior often introduce a low-rank structure in the covariance matrix.

Sparse GP Conclusions

- Exact GPs have an $\mathcal{O}(N^3)$ computational cost, making them feasible on small datasets with a few thousand instances only.
- Sparse GPs provide an approximate solution with a smaller computational cost that is $\mathcal{O}(NM^2)$ with $M \ll N$.
- The non-parametric property of GP is lost when using sparse approximations. They are no longer more flexible with more data.
- The methods that approximate the GP prior often introduce a low-rank structure in the covariance matrix.
- The best performing method seems to be the VFE method since it does not modify the prior.

Sparse GP Conclusions

- Exact GPs have an $\mathcal{O}(N^3)$ computational cost, making them feasible on small datasets with a few thousand instances only.
- Sparse GPs provide an approximate solution with a smaller computational cost that is $\mathcal{O}(NM^2)$ with $M \ll N$.
- The non-parametric property of GP is lost when using sparse approximations. They are no longer more flexible with more data.
- The methods that approximate the GP prior often introduce a low-rank structure in the covariance matrix.
- The best performing method seems to be the VFE method since it does not modify the prior.
- Some methods allow for stochastic optimization and mini-batch training that further reduce the cost to $\mathcal{O}(M^3)$.

References

- Williams, C., & Seeger, M. (2000). Using the Nyström method to speed up kernel machines. *Advances in neural information processing systems*, 13.
- Snelson, E., & Ghahramani, Z. (2005). Sparse Gaussian processes using pseudo-inputs. *Advances in neural information processing systems*, 18.
- Rahimi, A., & Recht, B. (2007). Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20.
- Naish-Guzman, A., & Holden, S. (2007). The generalized FITC approximation. *Advances in neural information processing systems*, 20.
- Hernández-Lobato, D., & Hernández-Lobato, J. M. (2016, May). Scalable Gaussian process classification via expectation propagation. In *Artificial Intelligence and Statistics* (pp. 168-176).
- Hensman, J., Fusi, N., & Lawrence, N. D. (2013). Gaussian processes for big data. *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*.
- Salimbeni, H., Eleftheriadis, S., & Hensman, J. (2018, March). Natural gradients in practice: Non-conjugate variational inference in Gaussian process models. In *International Conference on Artificial Intelligence and Statistics* (pp. 689-697).
- Wu, L., Miller, A., Anderson, L., Pleiss, G., Blei, D., & Cunningham, J. (2021). Hierarchical inducing point gaussian process for inter-domain observations. *International Conference on Artificial Intelligence and Statistics* (pp. 2926-2934).