

# Sparse Gaussian Processes: sgp

Carl Edward Rasmussen

April 1st, 2016

This document describes an implementation of two methods for sparse approximate inference in Gaussian Processes (GPs), the Fully Independent Training Conditionals (FITC) [1] and a Variational Free Energy approximation (VFE) [2].

The implementation assumes that the covariance function is squared exponential, and that the mean function is linear (for brevity, the mean function is not included in the documentation).

## 1 Log Marginal Likelihood

In the sparse approximation, the approximate log marginal likelihood is given by

$$\log(q(y|u)) = -\frac{1}{2}y^T(Q + G)^{-1}y - \frac{1}{2}\log|Q + G| - \frac{1}{2\sigma_n^2}\text{tr}(T) - \frac{n}{2}\log(2\pi), \quad (1)$$

where

$$Q_{f,f} = K_{f,u}K_{u,u}^{-1}K_{u,f}, \quad (2)$$

$$G_{\text{FITC}} = \text{diag}[K_{f,f} - Q_{f,f}] + \sigma_n^2 I, \quad \text{and} \quad G_{\text{VFE}} = \sigma_n^2 I, \quad (3)$$

$$T_{\text{FITC}} = 0, \quad \text{and} \quad T_{\text{VFE}} = K_{f,f} - Q_{f,f}. \quad (4)$$

Note that all elements of  $G$  are bounded below by  $\sigma_n^2$ , because  $K$  is positive definite. Rewriting  $Q_{f,f} = V^T V$ , where  $V = L^{-1}K_{u,f}$  and  $K_{u,u} = LL^T$ , and using the matrix inversion lemma, the approximate log marginal likelihood can be written as

$$\begin{aligned} \log(q(y|u)) &= -\frac{1}{2}y^T(G^{-1} - G^{-1}V^T A^{-1}VG^{-1})y - \frac{1}{2}\log|A| - \frac{1}{2}\log|G| - \frac{1}{2\sigma_n^2}\text{tr}(T) - \frac{n}{2}\log(2\pi) \\ &= -\frac{1}{2}y^T z - \frac{1}{2}\log|A| - \frac{1}{2}\log|G| - \frac{1}{2\sigma_n^2}\text{tr}(T) - \frac{n}{2}\log(2\pi) \end{aligned} \quad (5)$$

where  $A = I + VG^{-1}V^T$  and  $z = (Q + G)^{-1}y = (G^{-1} - G^{-1}V^T A^{-1}VG^{-1})y$ .

```
1 <nlml> 1)≡ (3a)
1 l = exp(hyp(e).l); s2 = exp(2*hyp(e).s); n2 = exp(2*hyp(e).n);
2 u = bsxfun(@rdivide, induce, 1'); % scaled inducing inputs
3 x = bsxfun(@rdivide, inputs, 1'); % scaled training inputs
4 Kuu = s2*(exp(-maha(u,u)/2) + ridge*eye(M));
5 Kuf = s2*exp(-maha(u,x)/2);
6 L = chol(Kuu)';
7 V = L\Kuf;
8 r = s2 - sum(V.*V,1)'; % diagonal residual Kff - Kfu Kuu^-1 Kuf
9 G = fitc*r + n2; iG = 1./G;
10 A = eye(M) + V*bsxfun(@times,iG,V');
11 J = chol(A)';
12 B = J\V;
13 z = iG.*y - (y'.*iG'*B'*B.*iG')';
14 nlml = nlml + y'*z/2 + sum(log(diag(J))) + sum(log(G))/2 ...
15 + vfe*sum(r)/n2/2 + N*log(2*pi)/2;
```

## 2 Derivatives

The derivative of the approximate log marginal likelihood wrt parameters  $\theta$

$$\frac{\partial \log(q(y|u))}{\partial \theta} = \frac{\partial \log q}{\partial Q} \frac{\partial Q}{\partial \theta} + \frac{\partial \log q}{\partial G} \frac{\partial G}{\partial \theta} + \frac{\partial \log q}{\partial T} \frac{\partial T}{\partial \theta}, \text{ where } q \triangleq q(y|u), \quad (6)$$

where  $\theta$  could be either the inducing inputs  $u$  or parameters of the covariance function (hyperparameters). We have

$$2 \frac{\partial \log q}{\partial Q} = \mathbf{z}\mathbf{z}^\top - \mathbf{G}^{-1} + \mathbf{G}^{-1}\mathbf{V}^\top \mathbf{A}^{-1}\mathbf{V}\mathbf{G}^{-1}, \text{ and } \frac{\partial \log q}{\partial G} = \text{diag}\left[\frac{\partial \log q}{\partial Q}\right]. \quad (7)$$

Note that we cannot compute  $\frac{\partial \log q}{\partial Q}$  itself, as the cost would be prohibitive ( $N^2M$ ).

### 2.1 Inducing inputs

We need the derivatives of  $Q$ ,  $G_{\text{FITC}}$  and  $T_{\text{VFE}}$  (ignoring the trivial  $G_{\text{VFE}}$  and  $T_{\text{FITC}}$ )

$$\frac{\partial Q}{\partial \mathbf{u}} = 2\mathbf{K}_{f,u}\mathbf{K}_{u,u}^{-1} \frac{\partial \mathbf{K}_{u,f}}{\partial \mathbf{u}} - \mathbf{K}_{f,u}\mathbf{K}^{-1} \frac{\partial \mathbf{K}_{u,u}}{\partial \mathbf{u}} \mathbf{K}^{-1} \mathbf{K}_{u,f} = 2 \text{sym}(\mathbf{R}^\top \mathbf{P}), \quad (8)$$

$$\frac{\partial G_{\text{FITC}}}{\partial \mathbf{u}} = -2 \text{diag}(\mathbf{R}^\top \mathbf{P}), \quad (9)$$

$$\frac{\partial T_{\text{VFE}}}{\partial \mathbf{u}} = -2 \text{tr}(\mathbf{R}^\top \mathbf{P}), \quad (10)$$

where  $\text{sym}(\mathbf{A}) = \frac{1}{2}(\mathbf{A} + \mathbf{A}^\top)$  and we have defined

$$\mathbf{R} = \mathbf{K}_{u,u}^{-1} \mathbf{K}_{u,f} \text{ and } \mathbf{P} = \frac{\partial \mathbf{K}_{u,f}}{\partial \mathbf{u}} - \frac{\partial \tilde{\mathbf{K}}_{u,u}}{\partial \mathbf{u}} \mathbf{R}, \quad (11)$$

and the derivative of  $\tilde{\mathbf{K}}_{u,u}$  denotes the derivative taken with respect to only the first argument of  $\mathbf{K}$  (as the derivative wrt the second argument is just the transpose of the derivative wrt the first argument). Of course, in an actual implementation, the product  $\mathbf{R}^\top \mathbf{P}$  in eq. (8) should never explicitly be computed (as this would cost  $N^2M$  operations).

2  $\langle \text{deriv } 2 \rangle \equiv$  (3a)

```

1 R = L'\V;
2 RiG = bsxfun(@times,R,iG');
3 RdQ = -R*z*z' + RiG - bsxfun(@times,RiG*B'*B,iG');
4 dG = z.^2 - iG + iG.^2.*sum(B.*B,1)';
5 RdQ2 = RdQ + bsxfun(@times, R, fitc*dG' - vfe/n2);
6 KW = Kuf.*RdQ2;
7 KWR = Kuu.*(RdQ2*R');
8 P = KW*x + bsxfun(@times, sum(KWR, 2) - sum(KW, 2), u) - KWR*u;
9 dnllml.induce = dnllml.induce + bsxfun(@rdivide, P, 1');
10 dnllml.hyp(e).l = -sum(P.*u,1) ...
11                 - sum((KW'*u - bsxfun(@times, sum(KW',2), x)).*x,1);
12 dnllml.hyp(e).n = -sum(dG)*n2 - vfe*sum(r)/n2;
13 dnllml.hyp(e).s = sum(sum(Kuf.*RdQ)) - fitc*r'*dG + vfe*sum(r)/n2;
14 dnllml.hyp(e).b = -sum(z);
15 dnllml.hyp(e).m = -inputs'*z;
```

## 2.2 Hyperparameters

For the log lengthscale hyperparameter

$$\frac{\partial Q}{\partial \log \ell} = \quad (12)$$

3a `<sgp.m 3a>`≡

```

1 function [nlml, dnlml] = sgp(p, inputs, target, style, test)
2
3 ridge = 1e-06; % relative jitter to make matrix better conditioned
4 switch style, case 'fitc', fitc = 1; vfe = 0; case 'vfe', vfe = 1; fitc = 0; end
5 induce = p.induce; hyp = p.hyp; % shorthand notation
6 [N, D] = size(inputs); M = size(induce,1); nlml = 0; dnlml.induce = zeros(M, D);
7
8 for e = 1:length(hyp)
9     y = target(:,e) - inputs*hyp(e).m - hyp(e).b;
10    <nlml 1>
11    if nargin == 5 % make predictions
12        <predict 3b>
13    elseif nargin == 2 % compute derivatives
14        <deriv 2>
15    end
16 end
```

## 3 Predictions for deterministic test inputs

For efficiency, predictions are made based on pre-computation of the quantities which don't depend on the test cases. These quantities are  $\beta$  and  $W$ . The Gaussian predictions

$$\mu = k(x^*, x)\beta, \quad \text{and} \quad \sigma^2 = k(x^*, x^*) - k(x^*, x)Wk(x, x^*). \quad (13)$$

The definition of  $\beta$  and  $W$  depend on the inference method being used. For the full GP we have

$$\beta = (K + \sigma_n^2 I)^{-1}y, \quad \text{and} \quad W = (K + \sigma_n^2 I)^{-1}. \quad (14)$$

For the sparse FITC and VFE method we have

$$\beta = (K_{u,u} + K_{u,f}G^{-1}K_{f,u})^{-1}K_{u,f}G^{-1}y, \quad \text{and} \quad W = K_{u,u}^{-1} - (K_{u,u} + K_{u,f}G^{-1}K_{f,u})^{-1}. \quad (15)$$

3b `<predict 3b>`≡ (3a)

```

1 beta = (y' .* iG' * B' / J / L)';
2 W = L' \ (eye(M) - eye(M) / J' / J) / L;
3 Ktu = s2 * exp(-maha(bsxfun(@rdivide, test, 1'), u) / 2);
4 nlml = Ktu * beta + test * hyp.m + hyp.b;
5 dnlml = s2 + n2 - sum(Ktu * W .* Ktu, 2);
```

## References

- [1] Edward Snelson and Zoubin Ghahramani. Sparse Gaussian processes using pseudo-inputs. In Y. Weiss, B. Schölkopf, and J. Platt, editors, Advances in Neural Information Processing Systems 18, pages 1257-1264. The MIT Press, Cambridge, MA, 2006.
- [2] Michalis K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. Twelfth International Conference on Artificial Intelligence and Statistics, (AISTATS), JMLR: W&CP 5, pp. 567-574, 2009.