



MAX-PLANCK-GESellschaft

Sparse GP's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Sparsity in Gaussian Processes Questions

Lehel Csató

Max-Planck Institute for Biological Cybernetics
Tübingen, Germany

Gaussian Processes Round Table¹



MAX-PLANCK-GESellschaft

Outline

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

- 1 Notations
 - Gaussian processes
 - Likelihoods
- 2 Approximations
 - Expectation propagation
 - Sparsification
 - Sparsity results
 - Comparisons
- 3 Examples
 - Predictive densities
 - Regression



Gaussian Processes

Sparse GP's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Gaussian processes – notations.

- For input locations $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, the associated random variables $\mathbf{f}_{\mathcal{X}} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^T$ are Gaussian:

$$p_0(\mathbf{f}_{\mathcal{X}}) = \mathbf{N}(\boldsymbol{\mu}_{\mathcal{X}}, \mathbf{K}_{\mathcal{X}})$$

- $\boldsymbol{\mu}_{\mathcal{X}}$ and $\mathbf{K}_{\mathcal{X}}$ are samples from the
 - Mean function $\mu(\mathbf{x}) = \langle f(\mathbf{x}) \rangle$
 - Covariance kernel function $K_0(\mathbf{x}, \mathbf{x}') = \langle f(\mathbf{x})f(\mathbf{x}') \rangle$
- $p_0(\mathbf{f}|\theta_1)$ denotes the *prior* process, θ_1 are parameters of the kernel function.



Data likelihood

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

- GP 's are *latent variables* in the inference process.
- Data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ is factorising with *likelihood function*:

$$P(\mathcal{D}|\mathbf{x}) = \prod_{n=1}^N P(y_n|f(\mathbf{x}_n), \theta_2)$$

θ_2 – parameters of the likelihood function

- The **posterior process** is computed using Bayes' rule

$$p_{\text{post}}(\mathbf{f}) = \frac{1}{Z} \int d\mathbf{f}_x P(\mathcal{D}|\mathbf{f}_x) p_0(\mathbf{f}_x, \mathbf{f})$$



Problems when computing

Sparse GP's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

$$p_{\text{post}}(\mathbf{f}) = \frac{1}{Z} \int d\mathbf{f}_x P(\mathcal{D}|\mathbf{f}_x) p_0(\mathbf{f}_x, \mathbf{f})$$

- Non-Gaussian likelihoods lead to non-Gaussian processes
 - $p_{\text{post}}(\mathbf{f})$ not analytically computable.
 - Cannot compute the normalising

$$Z = \int d\mathbf{f}_x P(\mathcal{D}|\mathbf{f}_x) p_0(\mathbf{f}_x)$$

- For Gaussian likelihoods

$$\mu_{\text{post}}(\mathbf{x}) = \mathbf{k}_N(\mathbf{x})^T \left(\sigma_o^2 \mathbf{I} + \mathbf{K}_{NN} \right)^{-1} \mathbf{y}$$

the matrix inversion becomes prohibitive.



Approximation steps

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

- Approximating the non-Gaussian posterior with a Gaussian one:
 - retains information about values *and uncertainties*.
 - retains non-Gaussianity when predicting:

$$p(y_* | \mathbf{x}_*, \mathcal{D}) = \int df_* P(y_* | f_*, \theta_2) p_{\text{post}}(f_*)$$

- Further approximation – with a **sparse support set**
 - keeps the information up to the second order.
 - speeds up the computation of $p_{\text{post}}(f_*)$.



“Best” approximation to the posterior

Sparse \mathcal{GP} 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Approximating the posterior process $p_{\text{post}}(\mathbf{f})$ with a Gaussian process $\hat{p}_{\text{post}}(\mathbf{f})$.

\iff minimising the Kullback-Leibler divergence:

$$d_{\text{KL}}(\hat{p}(\mathbf{f}) \| p_{\text{post}}(\mathbf{f})) = \int d\mathbf{f} p_{\text{post}}(\mathbf{f}) \log \frac{p_{\text{post}}(\mathbf{f})}{\hat{p}(\mathbf{f})}$$

The optimal \mathcal{GP} has the first and second moments of the non-Gaussian posterior:

$$\hat{\mu}(\mathbf{x}) = \mu_0(\mathbf{x}) + \sum_{n=1} \alpha_n K_0(\mathbf{x}, \mathbf{x}_n)$$

$$\hat{K}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') + \sum_{m,n=1} K_0(\mathbf{x}, \mathbf{x}_m) C_{mn} K_0(\mathbf{x}_n, \mathbf{x}')$$



Expectation propagation

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

- Iterating the following steps (ADATAP, Exp.Cons.):
 - For each n define the approximation that *excludes* n :
 $p_{\setminus n}(\mathbf{f})$
 - Build an approximation to the likelihood $P(y_n|f_n, \theta_2)$:

$$\frac{p_{\setminus n}(f_n) P(y_n|f_n, \theta_2)}{Z_n} \approx \mathcal{N}(f_n|\hat{\mu}_n, \hat{\sigma}_n) \stackrel{\text{def}}{=} \frac{p_{\setminus n}(f_n) \hat{t}(f_n|m_n, \lambda_n)}{\tilde{Z}_n}$$

- The variational step: finding parameters (m_n, λ_n) .
- Result – Gaussian approximation:

$$p_0(\mathbf{f}) \prod_n \frac{Z_n}{\tilde{Z}_n} \prod_n \hat{t}(f_n|m_n, \lambda_n)$$



The computational overload

Sparse \mathcal{GP} 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

The approximating \mathcal{GP} is “anchored” at \mathcal{X} :

$$\hat{K}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') - \sum K_0(\mathbf{x}, \mathbf{x}_m) \mathbf{C}_{mn} K_0(\mathbf{x}_n, \mathbf{x}')$$

- The computation time is **cubic** in data size.
- *Parameter* scale quadratically \iff over-parametrisation.

If data is “**structured**”, then there is a less redundant representation of *the same* approximation.



KL-optimal sparsification

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Sparse approximation step:

- The GP approximation $\hat{p}(\mathbf{f})$ is further reduced to a low-dimensional GP $p_{\mathcal{BV}}(\mathbf{f})$:

$$\hat{K}_{\mathcal{BV}}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') + \sum_{m, n \in \mathcal{BV}} K_0(\mathbf{x}, \mathbf{x}_m) \mathbf{C}_{mn} K_0(\mathbf{x}_n, \mathbf{x}')$$

- Kullback-Leibler divergence is used as minimiser:

$$p_{\mathcal{BV}}(\mathbf{f}) = \operatorname{argmin}_{p_d \in \mathcal{GP}_d} \text{KL}(p_d \| \hat{p}(\mathbf{f}))$$

\mathcal{GP}_d – GP 's with d locations.

- Optimisation is NP-complete, greedy sequential approach.



KL-optimal sparsification II.

Sparse \mathcal{GP} 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Assuming $|\mathcal{X}| = d + 1$

- Compute for each $n \in \mathcal{BV}_{d+1}$ (cheap approx.):

$$\text{KL}(p_n \parallel \hat{p}_{d+1}(\mathbf{f}))$$

Provides a measure of “how good” \mathbf{x}_n is.

- Remove the one with the minimum KL-loss.

Result

$$\hat{p}(\mathbf{f}) \propto p_0(\mathbf{f}) \prod_n \hat{t}(\pi_n \mathbf{f}_{\mathcal{BV}} | m_n, \lambda_n)$$

with

$$\pi_n \mathbf{f}_{\mathcal{BV}} = E[f_n | \mathbf{f}_{\mathcal{BV}}]_0$$



KL-optimisation results

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Reduced-rank GP using the **Basis Vectors – \mathcal{BV} set**:

$$p_{\mathcal{BV}}(\mathbf{f}) = (\mu_{\mathcal{BV}}, K_{\mathcal{BV}}, \mathcal{BV})$$

- Probabilistic framework;
- Selection of optimal \mathcal{BV} set via the “score” for $\mathbf{x} \in \mathcal{BV}$.
- Independent of the **noise/likelihood model**.

- **Sparse posterior variance shrunk.**
- Shrinkage – result of conditioning **but** is this good?
- A larger variance would probably be better.



MAX-PLANCK-GESellschaft

Comparisons

Sparse GP's

Lehel Csató

Notations

Gaussian processes
Likelihoods

Approximations

Expectation propagation
Sparsification
Sparsity results

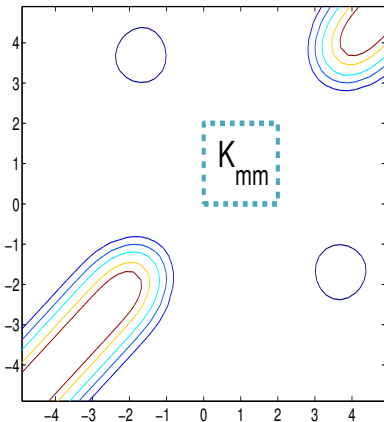
Comparisons

Examples

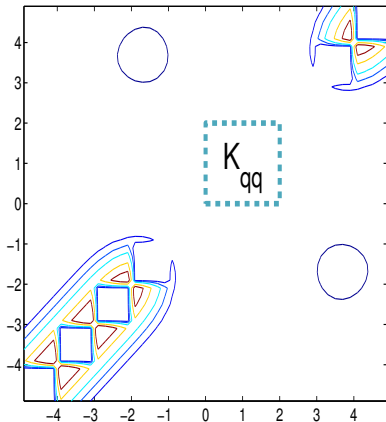
Predictive densities
Regression

Summary

GP approximations \iff viewed as approx to K_{NN}



Nystrom,
Sparse KL \equiv PPA



Bayesian Committee,
Snelson & Ghahramani



Sparse M.L.II.

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Optimising hyperparameters using

- 1 The log-Evidence is approximated as

$$\begin{aligned} \log \text{Ev} &= \sum_n \log Z_n - \sum \log \hat{Z}_n \\ &+ \log \int d\mathbf{f}_{B\mathcal{V}} p_0(\mathbf{f}_{B\mathcal{V}}) \prod_n \hat{t}(\pi_n \mathbf{f}_{B\mathcal{V}}) \end{aligned}$$

- 2 The upper bound to the log-evid (EM)

$$\log \text{Ev} = \int d\mathbf{f}_{B\mathcal{V}} p_{B\mathcal{V}}(\mathbf{f}) \log p_0(\mathbf{f}_{B\mathcal{V}}) P(\mathcal{D} | \Pi \mathbf{f}_{B\mathcal{V}})$$

The Evidence is better approximated using 1
but *the same test errors (class)*.



Inference Method

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

EM-EP:

Iterate the following steps

- Fix parameters \Rightarrow *approximate the posterior*
- Fix posterior \Rightarrow *new model parameters.*



Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

Efficiency:

- selecting a “good subset”.
- KL-based selection – depends on the (KL-)loss.
- Can be inefficient for different problems.
- Measures related to the *loss function* – see IVM



MAX-PLANCK-GESellschaft

Predictive densities for different models

Sparse GP's

Lehel Csató

- Gaussian likelihood involves no approximation

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

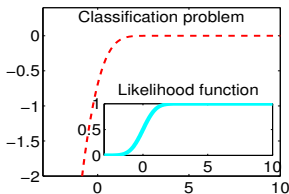
Examples

Predictive densities

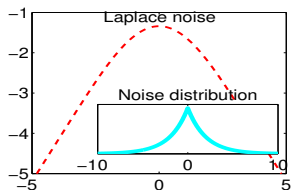
Regression

Summary

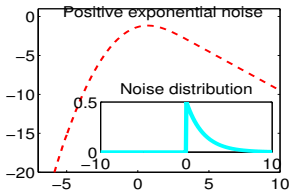
Classification



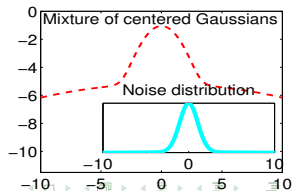
Laplace – robust



Pos. exponential



Mixture





Regression example

Sparse GP's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

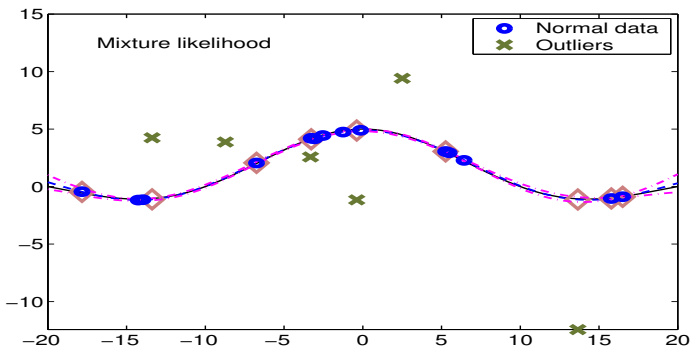
Summary

Detecting outliers using a mixture model:

$$P(y_n|f_n, \theta) = \pi N(y_n|f_n, \sigma_1^2) + (1 - \pi)N(y_n|f_n, \sigma_0^2)$$

Likelihood *not log-concave*.

π , σ_1^2 , and σ_0^2 estimated using MLII.





Concerns

Sparse GP's

Lehel Csató

Notations

Gaussian processes
Likelihoods

Approximations

Expectation propagation
Sparsification
Sparsity results
Comparisons

Examples

Predictive densities
Regression

Summary

- Subset selection important for GP's
- Sparsification based on KL-divergence might not be better
- Approximations – bounds on the Evidence – speed up computation.
- Selection criterion *should be based* on model **loss function or score**.



Summary

Sparse GP 's

Lehel Csató

Notations

Gaussian processes

Likelihoods

Approximations

Expectation propagation

Sparsification

Sparsity results

Comparisons

Examples

Predictive densities

Regression

Summary

- It is possible to infer the hyperparameters.
- Sparse approximation speeds up computation without significant loss.
- Outlook
 - Extension to two-level model specification.
 - Dynamical systems.

Software (matlab) and documentation available:

<http://www.tuebingen.mpg.de/~csato>