

Gaussian Processes

Perry Groot

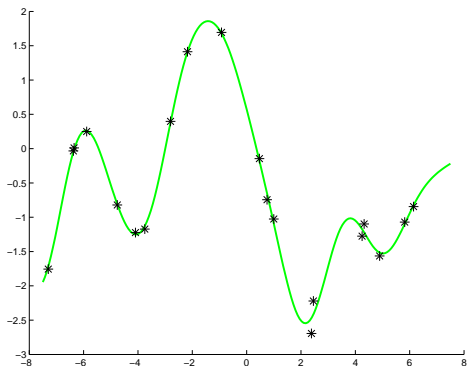
Radboud University Nijmegen

perry@cs.ru.nl

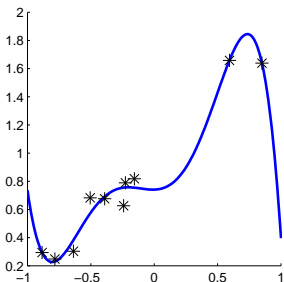
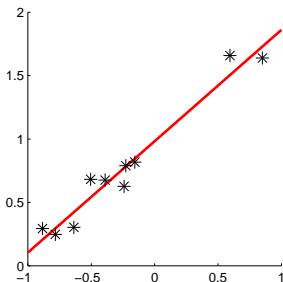
SIKS Course Computational Intelligence

11 & 12 March 2010

- Data $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$;
- Input space $\mathcal{X} \subseteq \mathbb{R}^d$; Output space $\mathcal{Y} \subseteq \mathbb{R}$
- Goal predict functional relation $f : \mathcal{X} \rightarrow \mathcal{Y}$



- *parametric* regression: $f(\mathbf{x}; \mathbf{w})$
 - Linear model: $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} = \sum_{j=0}^d w_j x_j$
 - Polynomial model: $f(x; \mathbf{w}) = \sum_{j=0}^M w_j x^j$
- Loss function: $\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$



GPs

Perry Groot

Regression

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

There are a couple of disadvantages:

- Lack of error bars on predictions
- Problem of overfitting

Overfitting can be avoided by using simpler models, but its predictive performance may be poor.

Bayes' rule to obtain a *posterior* distribution:

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2)p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X}, \sigma^2)}$$

predictive distribution

$$p(y_*|\mathbf{x}_*, \mathbf{y}, \mathbf{X}, \sigma^2) = \int p(y_*|\mathbf{x}_*, \mathbf{w}, \sigma^2)p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \sigma^2) d\mathbf{w}$$

- All parameters contribute to a prediction
- Good generalization performance and robust to overfitting
- Allows for error bars on predictions

GPs

Perry Groot

Regression

Gaussian processes

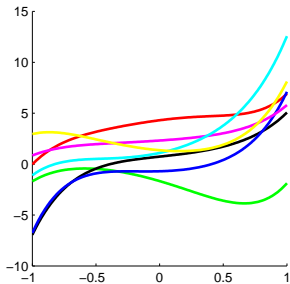
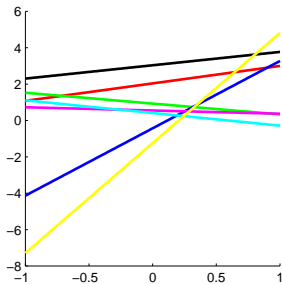
Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

Assuming a probability distribution over $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ leads to a probability distribution over functions $f(\cdot; \mathbf{w})$



GPs

Perry Groot

Regression

Gaussian processes

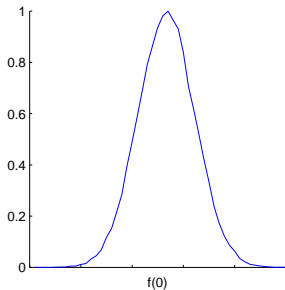
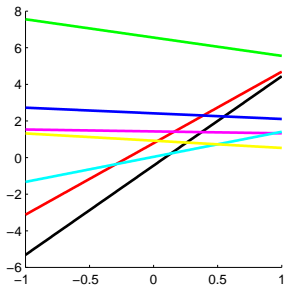
Posterior
Sampling
Model Selection

Classification

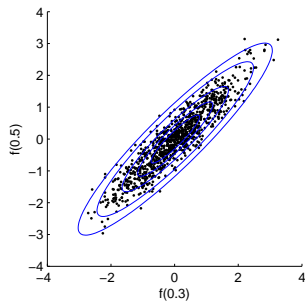
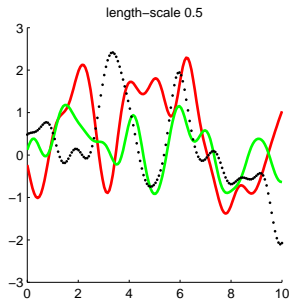
Applications

Preference Learning
Surrogate Modeling
Integration

Which leads to a distribution at each test point



Instead of taking a distribution over weights, we can also directly consider distributions over functions. We will consider the following model $y_i = f_i + \epsilon_i$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$



GPs

Perry Groot

Regression

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Applications

Preference Learning

Surrogate Modeling

Integration

A **Gaussian process** (GP) is collection of random variables $\{f_i\}$ with the property that the joint distribution of any finite subset has a joint Gaussian distribution.

A GP specifies a probability distribution over functions $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ and is fully specified by its mean function $m(\mathbf{x})$ and covariance (or kernel) function $k(\mathbf{x}, \mathbf{x}')$.

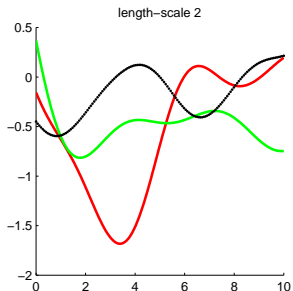
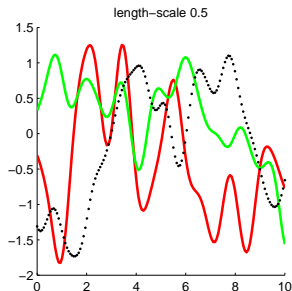
Typically $m(\mathbf{x}) = \mathbf{0}$, which gives

$$\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_l)\} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}) \text{ with } K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

Squared exponential (or Gaussian) covariance function:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\ell^2} \sum_{n=1}^N (x_n - x'_n)^2\right)$$

where ℓ is a length-scale parameter denoting how quickly the functions are to vary.



A priori, given data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ with $\mathbf{y} = f(\mathbf{X})$ and test points \mathbf{X}_* we have

$$\begin{bmatrix} f(\mathbf{X}) \\ f(\mathbf{X}_*) \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) & K(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

and after conditioning

$$f(\mathbf{X}_*) | \mathbf{X}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

with

$$\begin{aligned} \boldsymbol{\mu} &= K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{y} \\ \boldsymbol{\Sigma} &= K(\mathbf{X}_*, \mathbf{X}_*) - \underbrace{K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{X}_*)}_{\mathcal{O}(n^3)} \end{aligned}$$

GPs

Perry Groot

Regression

Gaussian processes

Posterior

Sampling

Model Selection

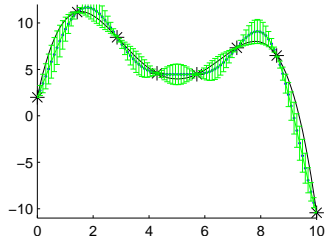
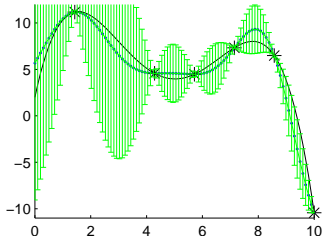
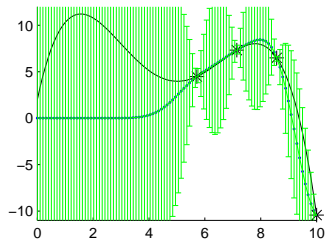
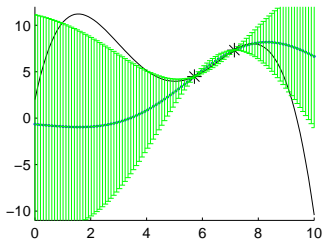
Classification

Applications

Preference Learning

Surrogate Modeling

Integration



GPs

Perry Groot

Regression

Gaussian processes

Posterior

Sampling

Model Selection

Classification

Applications

Preference Learning

Surrogate Modeling

Integration

How to sample functions from a $\mathcal{GP}(\mathbf{m}, \mathbf{K})$?

This can be done using the **Cholesky decomposition**, which is a lower triangular matrix \mathbf{L} such that $\mathbf{L}\mathbf{L}^T = \mathbf{K}$

$$\mathbf{L} = \text{chol}(\mathbf{K})^T;$$

$$\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I});$$

$$\mathbf{f} = \mathbf{m} + \mathbf{L}\mathbf{u}^T;$$

Then $\mathbb{E}[\mathbf{f}] = \mathbf{m} + \mathbf{L}\mathbb{E}[\mathbf{u}^T] = \mathbf{m}$ and

$$\text{var}(\mathbf{f}) = \text{var}(\mathbf{L}\mathbf{u}^T) = \mathbb{E}[\mathbf{L}\mathbf{u}^T\mathbf{u}\mathbf{L}^T] = \mathbf{L}\mathbb{E}[\mathbf{u}\mathbf{u}^T]\mathbf{L}^T = \mathbf{L}\mathbf{L}^T = \mathbf{K}$$

GPs

Perry Groot

Regression

Gaussian processes

Posterior

Sampling

Model Selection

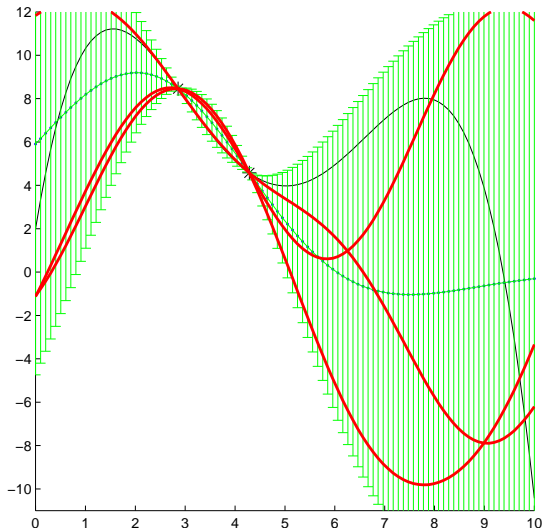
Classification

Applications

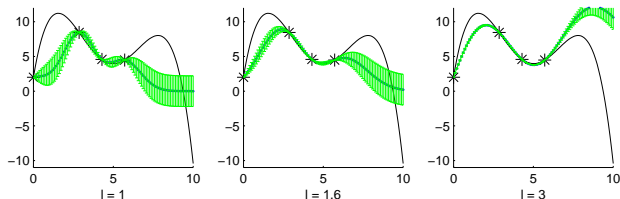
Preference Learning

Surrogate Modeling

Integration



The kernel function and likelihood may depend on additional parameters (**hyperparameters**) that need to be set



How to choose the best hyperparameters θ ?

In Bayesian model selection one can think of a hierarchical specification. At the bottom level, the posterior over parameters is

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \theta, \mathcal{H}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \theta, \mathcal{H})p(\mathbf{w}|\theta, \mathcal{H})}{p(\mathbf{y}|\mathbf{X}, \theta, \mathcal{H})}$$

where the evidence or marginal likelihood is

$$p(\mathbf{y}|\mathbf{X}, \theta, \mathcal{H}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathcal{H})p(\mathbf{w}|\theta, \mathcal{H}) d\mathbf{w}$$

Analogously, at the next level, the posterior over hyperparameters is

$$p(\theta|\mathbf{y}, \mathbf{X}, \mathcal{H}) = \frac{p(\mathbf{y}|\mathbf{X}, \theta, \mathcal{H})p(\theta|\mathcal{H})}{p(\mathbf{y}|\mathbf{X}, \mathcal{H})}$$

- This hierarchical specification can be continued analogously to obtain a posterior over models \mathcal{H} .
- Depending on the model, however, some integrals may be intractable and approximations are needed.
- In a type II maximum likelihood approximation we *maximize the marginal likelihood*.

$$\log p(\mathbf{y}|\mathbf{X}, \theta) = -\frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \mathbf{y}^T \mathbf{K} \mathbf{y} - \frac{n}{2} \log 2\pi$$

GPs

Perry Groot

Regression

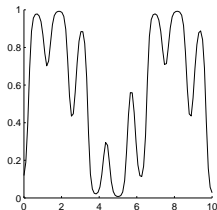
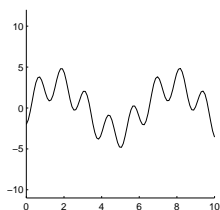
Gaussian
processesPosterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

- GPs can also be used for classification, but computations are intractable (needs approximations).
- The idea is to squash a regression function in the domain $(-\infty, \infty)$ to the domain $[0, 1]$
 - Logistic regression: $\lambda(\mathbf{x}^T \mathbf{w})$ with $\lambda(z) = \frac{1}{1 + \exp(-z)}$
 - Probit regression: $\Phi(z) = \int_{-\infty}^z \mathcal{N}(x|0, 1) dx$



GPs

Perry Groot

Regression

Gaussian processes

Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

- Regression / Classification
- Clustering
- Ordinal Regression
- Preference Learning
- Ranking
- Surrogate modeling / Global Optimization
- Relational Learning
- Reinforcement Learning
- Visualization of high dimensional data
- Nonrigid Shape Recovery
- Evaluating Integrals
- ...

Problem: Given a data set of M pairwise preferences (i.e., a set of pairs $(\mathbf{x}_1, \mathbf{x}_2)$ and whether $\mathbf{x}_1 \succ \mathbf{x}_2$ or $\mathbf{x}_1 \prec \mathbf{x}_2$ holds)

$$\mathcal{D} = \{(\mathbf{x}_{m_1}, \mathbf{x}_{m_2}, d_m) \mid 1 \leq m \leq M, d_m \in \{-1, 1\}\}$$

predict for new instances \mathbf{x}, \mathbf{y} which one is preferred.

Idea: Assume a latent (utility) function f over instances that preserves user preferences, i.e., basically $f(\mathbf{x}_1) > f(\mathbf{x}_2)$ when $\mathbf{x}_1 \succ \mathbf{x}_2$.

Bayesian framework

$$p(\mathbf{f}|\mathcal{D}, \mathcal{H}) = \frac{p(\mathbf{f}|\mathcal{H})p(\mathcal{D}|\mathbf{f}, \mathcal{H})}{p(\mathcal{D}|\mathcal{H})}$$

with a likelihood function, for $b \in \mathbb{R}$, $\delta_1, \delta_2 \sim \mathcal{N}(0, \sigma^2)$,

$$\begin{aligned} p(\mathbf{x}_1 \succ \mathbf{x}_2 | f(\mathbf{x}_1), f(\mathbf{x}_2)) &= p(f(\mathbf{x}_1) + \delta_1 > f(\mathbf{x}_2) + b + \delta_2) \\ &= \Phi(z) \end{aligned}$$

with

$$z = \frac{d(f(\mathbf{x}_1) - f(\mathbf{x}_2) - b)}{\sqrt{2}\sigma}$$

GPs

Perry Groot

Regression

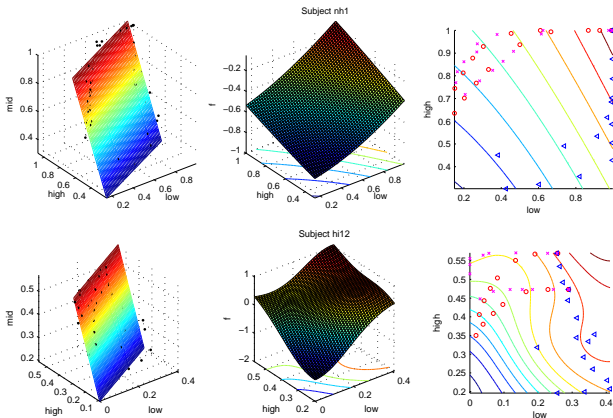
Gaussian processes

Posterior Sampling
Model Selection

Classification

Applications
Preference Learning
Surrogate Modeling
Integration

Applied to 14 normal-hearing and 18 hearing-impaired subjects. Obtained significant improvement for predicting preferences of hearing-impaired subjects.



GPs

Perry Groot

Regression

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

Complex (physical) systems can be studied nowadays by computer simulations, but often need long running times.

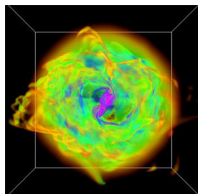
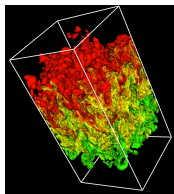
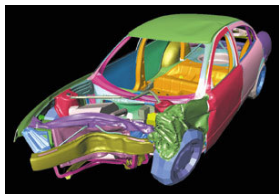


Figure: 1. Car collision; 2. Turbulent-mixing dynamics of a supernova; 3. Gas cloud collapsing inwards to form a star.

GPs

Perry Groot

Regression

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Applications

Preference Learning

Surrogate Modeling

Integration

Idea: replace costly simulation model by a fast Gaussian process surrogate model.

Choose function evaluations in a "smart way" (e.g., reducing overall variance) to obtain a good model fit.

Sometimes, however, we are not interested in a good global model, but only in one specific point (e.g., the best parameter setting).

$$\tilde{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$$

Let $f_{\max} = \max\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)\}$ be the best value so far. The **improvement** at a new point $y = f(\mathbf{x})$ is defined as

$$I(\mathbf{x}) = \max\{0, f(\mathbf{x}) - f_{\max}\}$$

Using the GP prediction $y = f(\mathbf{x}) \sim \mathcal{N}(m, s^2)$ we obtain the **Expected Improvement (EI)**:

$$E(I) = \begin{cases} (m - f_{\max})(1 - \Phi(d)) + s\phi(d) & s > 0 \\ 0 & s = 0 \end{cases}$$

with $d = (f_{\max} - m)/s$ and where $\Phi(\cdot)$ and $\phi(\cdot)$ denote the cdf and pdf of the standard normal distribution.

GPs

Perry Groot

Regression

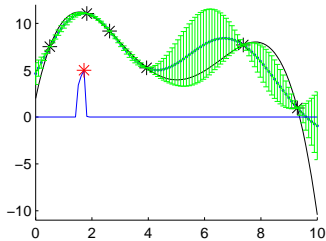
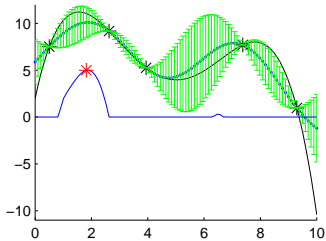
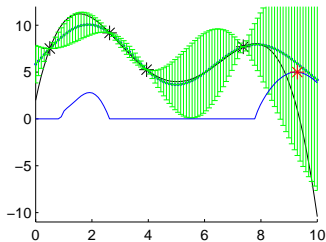
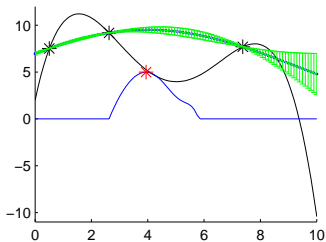
Gaussian processes

Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration



Advantages and disadvantages EI:

- EI allows for exploration and exploitation
- Expected Improvement can be sampled fast
- Expected Improvement often converges to a local optimum

Generalized Improvement (g positive integer):

$$I^g(\mathbf{x}) = \begin{cases} (f(\mathbf{x}) - f_{\max})^g & f(\mathbf{x}) > f_{\max} \\ 0 & \text{otherwise} \end{cases}$$

Larger g results in more exploration.

GPs

Perry Groot

Regression

Gaussian
processesPosterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

Consider the following integral

$$F = \int_{\mathbf{x}} f(\mathbf{x})p(\mathbf{x}) d\mathbf{x}$$

where p is a known distribution over the inputs X .

For example,

- f could be a computer simulation needing several hours of computation to evaluate f in a single point
- $p(\mathbf{x})$ is the posterior distribution and $f(\mathbf{x})$ predictions made by the model with parameters \mathbf{x} ,
- $p(\mathbf{x})$ is the parameter prior and $f(\mathbf{x}) = p(y|\mathbf{x})$ the likelihood (i.e., integral computes the evidence)

GPs

Perry Groot

Regression

Gaussian
processesPosterior
Sampling
Model Selection

Classification

Applications
Preference Learning
Surrogate Modeling
Integration

Monte Carlo makes the approximation

$$F \simeq \frac{1}{T} \sum_{t=1}^T f(\mathbf{x}^{(t)})$$

with $\mathbf{x}^{(t)}$ random draws from $p(\mathbf{x})$. Disadvantages are (c.f. [O'Hagan, 1987] 'Monte Carlo is Fundamentally Unsound'):

- MC is a frequentist approach
- MC can use an irrelevant *importance sampling distribution* $q(\mathbf{x})$ when sampling is hard from $p(\mathbf{x})$
- MC ignores the values $\mathbf{x}^{(t)}$

GPs

Perry Groot

Regression

Gaussian
processesPosterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

We can think of F as being **random** as we are uncertain about $f(\mathbf{x})$ because we cannot afford to compute $f(\mathbf{x})$ at every location.

The integral is then a **Bayesian inference problem**:

- put a prior on f ,
- for observations, evaluate f in a number of points
- combine the prior and observations into a posterior distribution over f (which implies a distribution over F)

When the prior f and posterior $f|\mathcal{D}$ are GPs, the distribution of F is Gaussian, $F \sim \mathcal{N}(\bar{F}, \text{cov}(F))$, and is fully characterized by its mean and variance

$$\begin{aligned} \bar{F} &= \int_{\mathbf{x}} \bar{f}_{\mathcal{D}}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ \text{cov}(F) &= \int_{\mathbf{x}} \int_{\mathbf{x}'} \text{cov}(f_{\mathcal{D}}(\mathbf{x}), f_{\mathcal{D}}(\mathbf{x}')) p(\mathbf{x}) p(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \end{aligned} \quad (1)$$

with $\bar{f}_{\mathcal{D}}$ and $\text{cov}(f_{\mathcal{D}}(\mathbf{x}), f_{\mathcal{D}}(\mathbf{x}'))$ the posterior mean and posterior variance, respectively.

Sometimes the problem can be reduced to products of one dimensional integrals and/or some analytic expression, e.g.,

$$p(\mathbf{x}) \sim \mathcal{N}(\mathbf{b}, \mathbf{B})$$

$$k(\mathbf{x}, \mathbf{x}') = w_0 \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \mathbf{A}^{-1}(\mathbf{x} - \mathbf{x}')\right)$$

with $\mathbf{A} = \text{diag}(w_1^2, \dots, w_N^2)$. Then

$$\bar{F} = \mathbf{zK}^{-1} Y, \quad \text{cov}(F) = k_c - \mathbf{zK}^{-1} \mathbf{z}^T$$

with

$$k_c = w_0 |2\mathbf{A}^{-1} \mathbf{B} + \mathbf{I}|^{-1/2}$$

$$z_l = w_0 |\mathbf{A}^{-1} \mathbf{B} + \mathbf{I}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x}_l - \mathbf{b})^T (\mathbf{A} + \mathbf{B})^{-1} (\mathbf{x}_l - \mathbf{b})\right)$$

GPs

Perry Groot

Regression

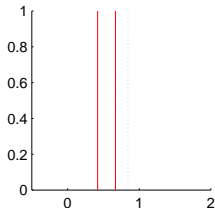
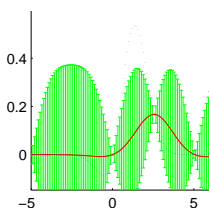
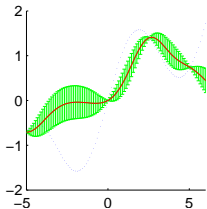
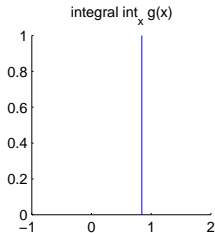
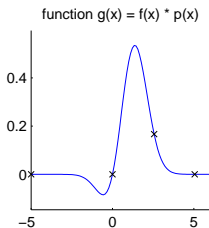
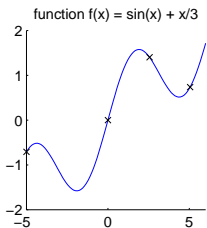
Gaussian processes

Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration



GPs

Perry Groot

Regression

Gaussian processes

Posterior Sampling
Model Selection

Classification

Applications
Preference Learning
Surrogate Modeling
Integration

Suppose we want to introduce a new cake mix into the consumer market which is robust to an inaccurate setting of oven temperature and baking time.

3 Control variables: The amount of flour (F), the amount of sugar (S), and the amount of egg powder (E).

2 Noise variables: Oven temperature (T) and baking time (t).



GPs

Perry Groot

Regression

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Applications

Preference Learning
Surrogate Modeling
Integration

- Gaussian Processes for Machine Learning,
C.E. Rasmussen and K.I. Williams
- The Gaussian Process Web Site
<http://www.gaussianprocess.org/>