

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Introduction to Bayesian Inference and Gaussian Processes

Perry Groot

Radboud University Nijmegen

`perry@cs.ru.nl`

Computational Intelligence and Learning Doctoral School
Université catholique de Louvain, Louvain-la-Neuve
3 Feb 2014

Linear Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- 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 is to:
 - *Learn* functional relationship between \mathcal{X} and \mathcal{Y}

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

- *Predict* unknown target values given new input values

Linear Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- How do we learn a functional relationship from a finite number of observations?
- Given a model, how do we determine the predictive quality of the model?

Linear Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Linear regression model:

$$f(x; w_0, w_1) = w_1 x + w_0$$

- Values for free parameters w_0, w_1 need to be defined given the observed data

Loss function

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Identification of model parameters can be done by a **loss function** that defines the miss-match between the output of the model and observed target values.
- Mean Squared error loss

$$\frac{1}{N} \sum_{n=1}^N (y_n - f(x_n; w_0, w_1))^2$$

- Mean Absolute error loss

$$\frac{1}{N} \sum_{n=1}^N |y_n - f(x_n; w_0, w_1)|$$

Loss function

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Let $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$, $\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$

- Then $f(\mathbf{X}; \mathbf{w}) = \mathbf{X}\mathbf{w}$

- Mean Squared Error (MSE) loss

$$\frac{1}{N}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$$

- How do we minimize the MSE?

Loss function

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- $\hat{\mathbf{w}}$ minimizes the MSE when the gradient of all its partial derivatives is zero

$$\begin{aligned}\frac{\partial \text{MSE}}{\partial \mathbf{w}} &= \left[\frac{\partial \text{MSE}}{\partial w_0} \quad \frac{\partial \text{MSE}}{\partial w_1} \right] = \left[\frac{2}{N} \sum_{n=1}^N (f(x_n; w_0, w_1) - y_n) \quad \frac{2}{N} \sum_{n=1}^N (f(x_n; w_0, w_1) - y_n) x_n \right] \\ &= -\frac{2}{N} \mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{w}) = \mathbf{0} \\ &\Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

- The Hessian $\frac{\partial^2 \text{MSE}}{\partial \mathbf{w} \partial \mathbf{w}^T} = \frac{2}{N} \mathbf{X}^T \mathbf{X}$ is positive definite when $\mathbf{X}^T \mathbf{X}$ is invertible, and therefore $\hat{\mathbf{w}}$ is a minimum

Polynomial Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- We could also assume a polynomial form

$$f(x; \mathbf{w}) = w_k x^k + \dots + w_2 x^2 + w_1 x + w_0 = \sum_{i=0}^K w_i x^i$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^N \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_N & x_N^2 & \dots & x_N^N \end{bmatrix}$$

- Model linear in parameters
- Non-linear model because of a non-linear transformation of the inputs

Linear and Polynomial Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

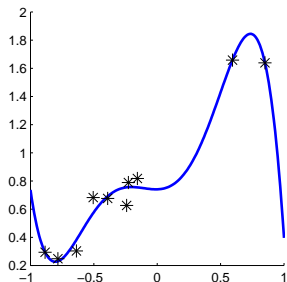
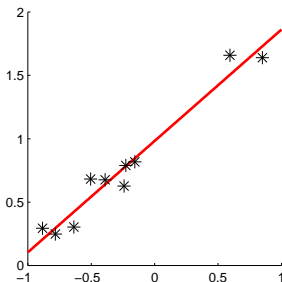
Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- parametric regression model: $f(\mathbf{x}; \mathbf{w})$

- Linear model: $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} = \sum_{j=0}^d w_j x_j$

- Polynomial model: $f(x; \mathbf{w}) = \sum_{j=0}^K w_j x^j$

- Loss function: $\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$



Supervised Learning: Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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.

Probability Theory

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Probabilities provide a means to represent uncertainty, e.g., ‘the probability of rolling a 6 with a die is $1/6$ ’.
- Two views: frequentist and Bayesian
 - Frequentist: frequency in long run of experiments
 - Bayesian: a degree of belief

Probability Theory

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- probabilities are non-negative: $p(x) \geq 0$,
- probabilities normalize: $\sum_x p(x) = 1$,
- sum rule: $p(x) = \sum_y p(x, y)$,
- product rule: $p(x, y) = p(x|y)p(y) = p(y|x)p(x)$,
- joint probability distribution: $p(x, y)$,
- conditional probability: $p(x|y) = p(x, y)/p(y)$,
- Bayes rule: $p(y|x) = p(x|y)p(y)/p(x)$.

Probability Densities

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- A continuous variable X has a probability density function (pdf) f_X when

$$p(x \in (a, b)) = \int_a^b f_X(x) dx$$

- and has cumulative distribution function F_X when

$$F_X(x) = \int_{-\infty}^x f_X(u) du$$

Probability Densities

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

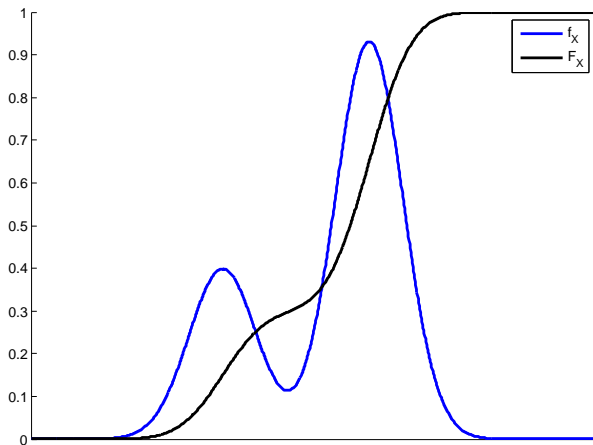
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Gaussian Distribution

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

■ probability density functions

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

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

■ cumulative density function

$$\Phi(x; \mu; \sigma^2) = \Phi((x - \mu)/\sigma) = \int_{-\infty}^x \phi(z; \mu, \sigma^2) dz$$

Gaussian Distribution

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Let x and y have a joint normal distribution

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix} \right)$$

Then the marginal distribution of x is

$$x \sim \mathcal{N}(\mu_x, A)$$

and the conditional distribution of x given y is

$$x|y \sim \mathcal{N}(\mu_x + CB^{-1}(y - \mu_y), A - CB^{-1}C^T)$$

Expectation and Variance

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

■ In general

$$\mathbb{E}[f] = \sum_x p(x)f(x) \quad \text{or} \quad \mathbb{E}[f] = \int p(x)f(x)dx$$

$$\text{var}[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] = \mathbb{E}[f(x)^2] - [\mathbb{E}[f(x)]]^2$$

■ Notation: $\langle f \rangle_q = \mathbb{E}_q[f] = \int f q(f)df$

Probabilistic Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Express uncertainty over the target values using a probability distribution

$$y = f(x; \mathbf{w}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- Interested in how likely are the observed outputs given the inputs and model parameters
- Likelihood of an observation is the conditional probability $p(y|x, \mathbf{w}, \sigma)$
- Data likelihood (assuming independent measurements) is given by the **likelihood function**

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma) = \prod_{n=1}^N p(y_n|x_n, \mathbf{w}, \sigma) = \prod_{n=1}^N \mathcal{N}(y_n; f(x_n; \mathbf{w}), \sigma)$$

Probabilistic Regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma) = \prod_{n=1}^N \mathcal{N}(y_n; f(x_n; \mathbf{w}), \sigma)$$

- Data can be made more likely under the model by optimizing the parameters
- Easier to use log-likelihood $\mathcal{L} = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma)$

$$-\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(x_n; \mathbf{w}))^2$$

Maximum Likelihood

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Log-likelihood is maximized by setting the partial derivatives to 0

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \mathbf{w}) = \mathbf{0}$$

- Hessian is strictly negative implying a maximum

$$\frac{\partial^2 \mathcal{L}}{\partial \mathbf{w} \partial \mathbf{w}^T} = -\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X}$$

- Maximum-likelihood (ML) solution

$$\hat{\mathbf{w}}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Maximum Likelihood

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- The ML estimate $\hat{\sigma}_{ML}$ can be obtained analogously
- Since we have a probabilistic model, new predictions are expressed in terms of a **predictive distribution** instead of a point estimate

$$p(y|x, \hat{\mathbf{w}}_{ML}, \hat{\sigma}_{ML}) = \mathcal{N}(y; f(x; \hat{\mathbf{w}}_{ML}), \hat{\sigma}_{ML}^2)$$

Maximum a Posteriori

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Could go one step further by introducing a distribution over parameters

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \alpha^{-1} \mathbf{I})$$

where the parameters controlling the distributions of parameters are called **hyperparameters**

- Using Bayes' rule $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \alpha, \sigma) \propto p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma)p(\mathbf{w}|\alpha)$ the most probable value $\hat{\mathbf{w}}_{MAP}$ can be obtained using similar strategies

Bayesian regression

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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

Weightspace view

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

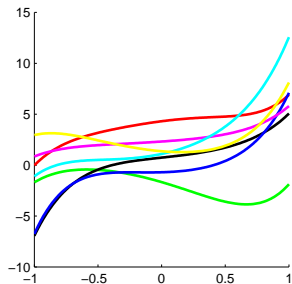
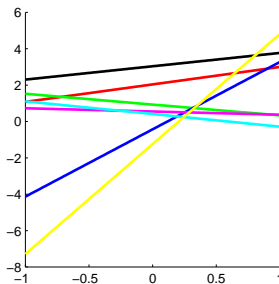
Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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



Weightspace view

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

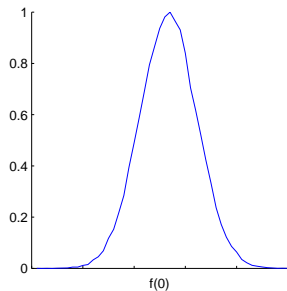
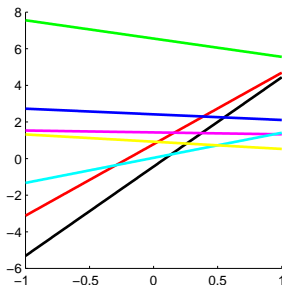
Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Which leads to a distribution at each test point



Functionspace view

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

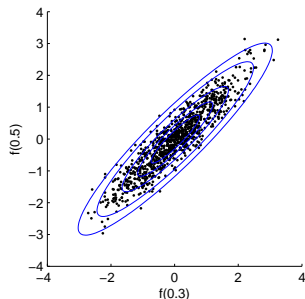
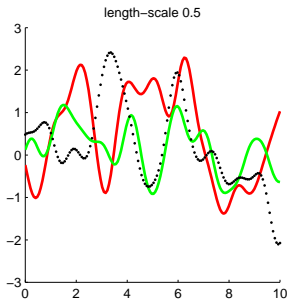
Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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 $\epsilon \sim \mathcal{N}(0, \sigma^2)$



Gaussian Processes

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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

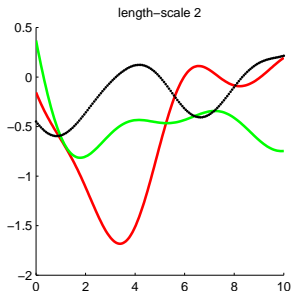
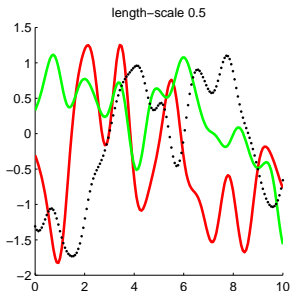
$$\{f(\mathbf{x}_1), \dots, f(\mathbf{x}_I)\} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}) \text{ with } K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

Gaussian Processes - Covariance function

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.



Gaussian Distribution

Perry Groot

Regression

Probabilistic
Inference

**Gaussian
processes**

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Let's consider a Gaussian with a particular distribution
- Generate a single sample from this 25 dimensional distribution $f = [f_1, f_2, \dots, f_{25}]$
- plot these samples against their indexes

Gaussian Distribution Sample

Perry Groot

Regression

Probabilistic
Inference

**Gaussian
processes**

Posterior
Sampling
Model Selection

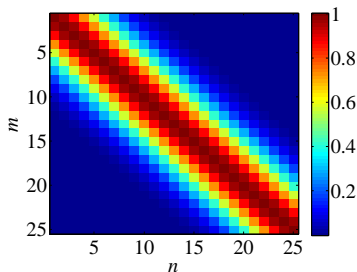
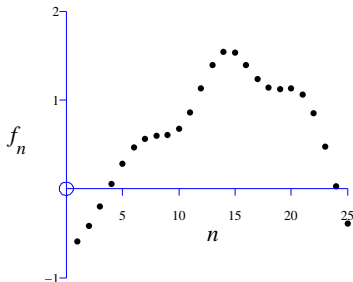
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



Left: 25 dimensional correlated random variable plotted against index. Right: colormap showing correlations between dimensions.

Covariance function

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Covariance matrix shows correlation between points f_i and f_j if i is near to j .
- Less correlation if i is distant from j .
- Ordering of points means that the function appears smooth.
- We will focus on the distribution of two points.

Prediction of f_2 from f_1

Perry Groot

Regression

Probabilistic
Inference

**Gaussian
processes**

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

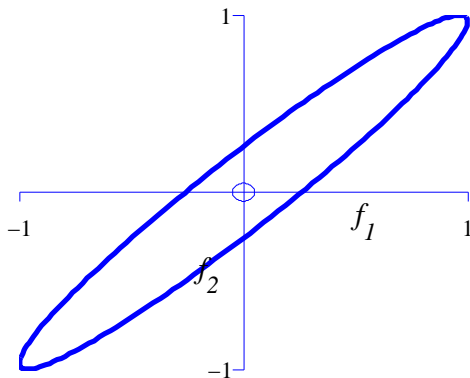


Figure: Covariance for $\begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$ is $K_{12} = \begin{bmatrix} 1 & 0.966 \\ 0.966 & 1 \end{bmatrix}$.

Prediction of f_2 from f_1

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

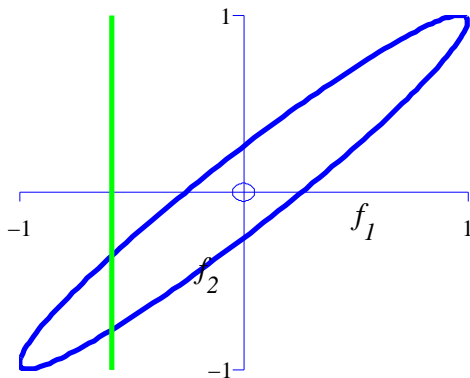


Figure: Covariance for $\begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$ is $K_{12} = \begin{bmatrix} 1 & 0.966 \\ 0.966 & 1 \end{bmatrix}$.

Prediction of f_2 from f_1

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

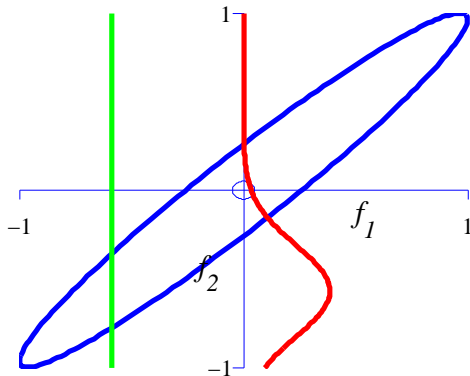


Figure: Covariance for $\begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$ is $K_{12} = \begin{bmatrix} 1 & 0.966 \\ 0.966 & 1 \end{bmatrix}$.

Prediction of f_5 from f_1

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

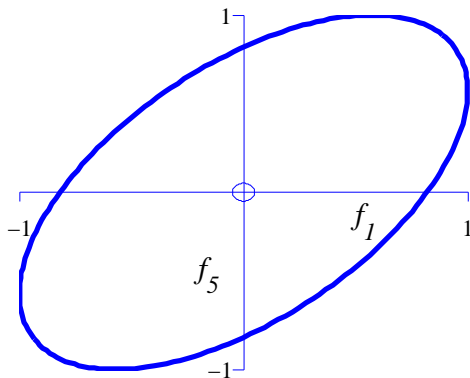


Figure: Covariance for $\begin{bmatrix} f_1 \\ f_5 \end{bmatrix}$ is $K_{15} = \begin{bmatrix} 1 & 0.574 \\ 0.574 & 1 \end{bmatrix}$.

Prediction of f_5 from f_1

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

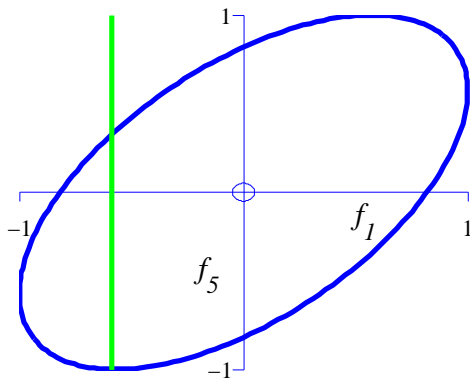


Figure: Covariance for $\begin{bmatrix} f_1 \\ f_5 \end{bmatrix}$ is $K_{15} = \begin{bmatrix} 1 & 0.574 \\ 0.574 & 1 \end{bmatrix}$.

Prediction of f_5 from f_1

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

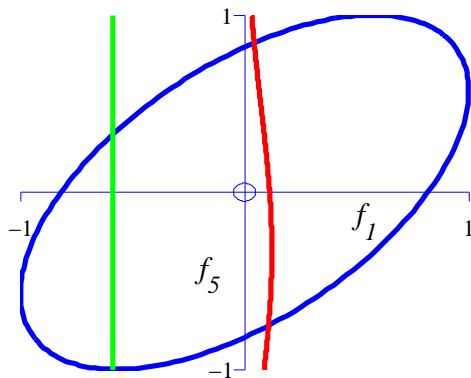


Figure: Covariance for $\begin{bmatrix} f_1 \\ f_5 \end{bmatrix}$ is $K_{15} = \begin{bmatrix} 1 & 0.574 \\ 0.574 & 1 \end{bmatrix}$.

Gaussian Processes - Posterior process

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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

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

Gaussian Processes - 1D demo

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

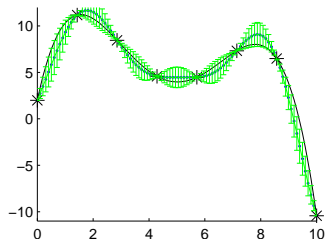
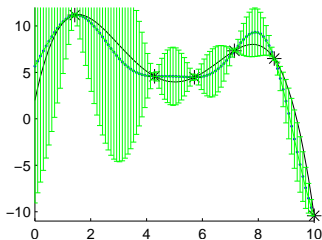
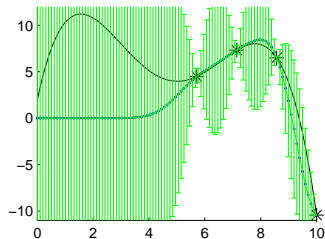
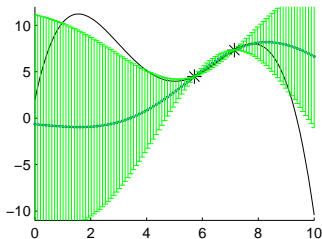
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Gaussian Processes - Sampling

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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

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

Gaussian Processes - Sampling

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

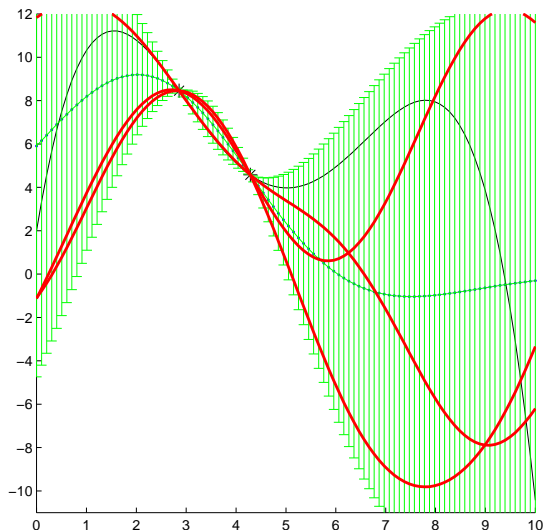
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Model Selection: Hyperparameters

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

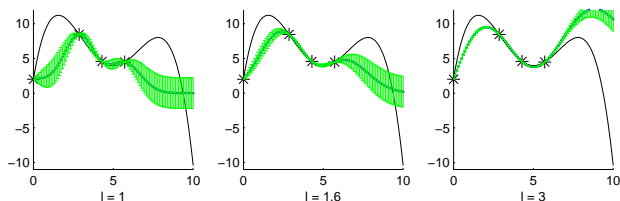
Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

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



How to choose the best hyperparameters θ ?

Model Selection: Marginal Likelihood

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

For learning kernel parameters we typically optimize the *maximize the marginal likelihood*. For regression:

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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

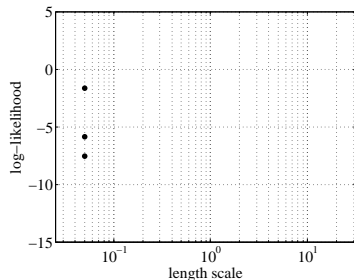
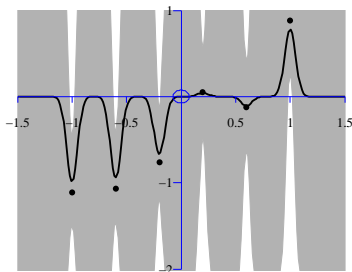
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

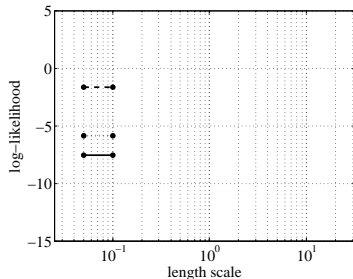
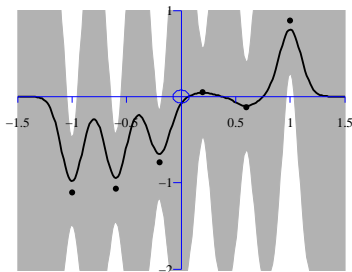
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

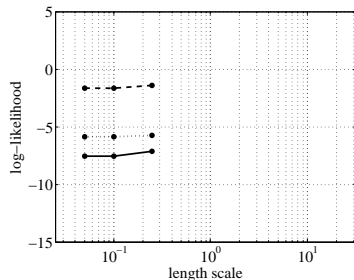
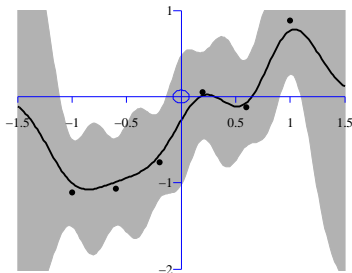
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

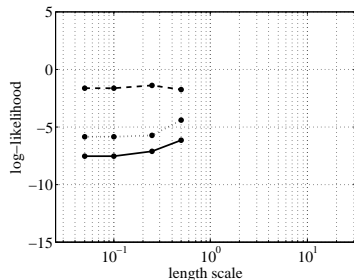
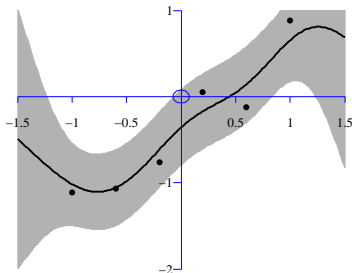
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

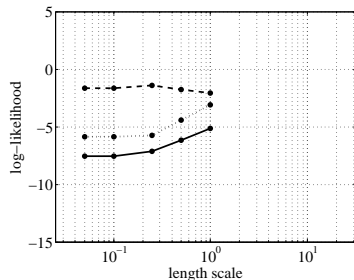
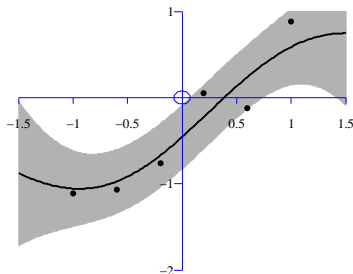
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

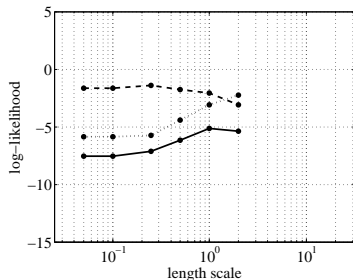
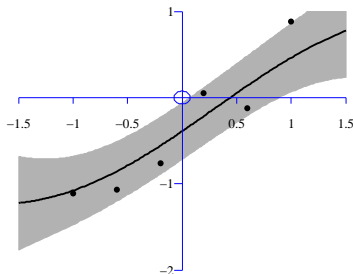
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

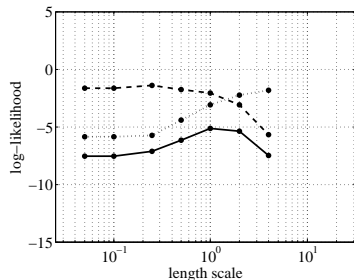
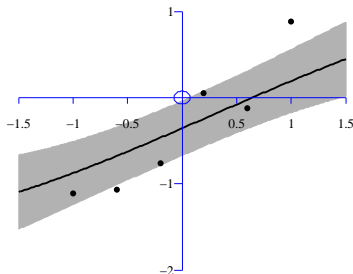
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

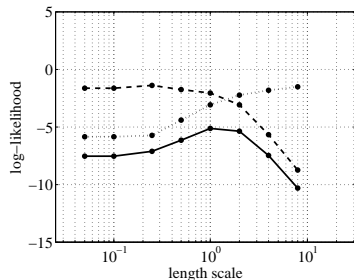
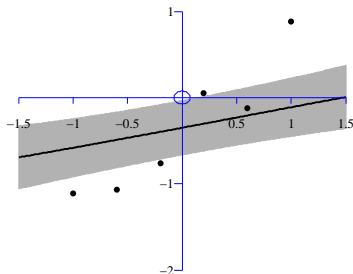
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Model Selection: Example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

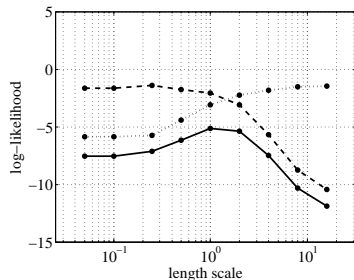
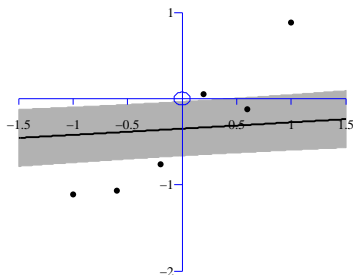
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



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

Classification

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

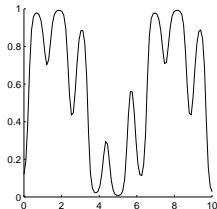
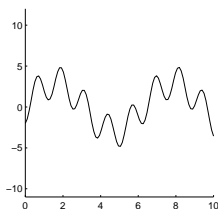
Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- 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: $\Phi(z) = \int_{-\infty}^z \mathcal{N}(x|0, 1) dx$



Laplace Approximation

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Approximates the posterior at the maximum a posteriori (MAP) estimate of the latent functions: \hat{f}

$$\begin{aligned}\hat{f} &= \operatorname{argmax}_f p(f|\mathcal{D}) \\ &= \operatorname{argmax}_f \frac{p(\mathcal{D}|f)p(f)}{p(\mathcal{D})} \\ &= \operatorname{argmax}_f \log \left(\frac{p(\mathcal{D}|f)p(f)}{p(\mathcal{D})} \right) \\ &= \operatorname{argmax}_f \log p(\mathcal{D}|f) + \log p(f) \\ &= \operatorname{argmax}_f \Psi(f)\end{aligned}$$

Laplace Approximation

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

$$\begin{aligned}\Psi(f) &= \log p(\mathcal{D}|f) + \log p(f) \\ &= \log p(\mathcal{D}|f) - \frac{1}{2}f^T K^{-1}f - \frac{1}{2}\log |K| - \frac{n}{2}\log(2\pi)\end{aligned}$$

$$\begin{aligned}\Psi(f) &\simeq \Psi(\hat{f}) + \frac{1}{2}(f - \hat{f})^T \nabla \nabla \Phi(\hat{f})(f - \hat{f}) \\ &\simeq \Psi(\hat{f}) - \frac{1}{2}(f - \hat{f})^T [K^{-1} + W](f - \hat{f})\end{aligned}$$

$$\begin{aligned}p(f|\mathcal{D}) &\propto p(\mathcal{D}|f)p(f) = \exp(\Psi(f)) \\ &\simeq \exp\left(\Psi(\hat{f}) - \frac{1}{2}(f - \hat{f})^T [K^{-1} + W](f - \hat{f})\right) \\ &\propto \mathcal{N}(\hat{f}, (K^{-1} + W)^{-1})\end{aligned}$$

Expectation Propagation

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- The posterior $p(f|\mathcal{D}) = \frac{p(f)p(\mathcal{D}|f)}{p(\mathcal{D})}$ is **intractable**
- EP approximates the likelihood by a Gaussian distribution making the posterior tractable
- Local likelihood approximations

$$p(y_i|f_i) \simeq t_i(f_i|\tilde{Z}_i, \tilde{\mu}_i, \tilde{\sigma}_i^2) = \tilde{Z}_i \mathcal{N}(f_i|\tilde{\mu}_i, \tilde{\sigma}_i^2)$$

- Approximation is iteratively updated
- In the Gaussian case the update step turns out to be the same as **moment matching**

Variational Approximation

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Given observations y and latent variables f , in variational inference we have the following fundamental relation

$$\mathcal{F}(q(f)) = \log p(y) - \text{KL}(q(f) \| p(f|y))$$

where

$$\mathcal{F}(q(f)) = \int q(f) \log \frac{p(f, y)}{q(f)} df = \left\langle \log \frac{p(f, y)}{q(f)} \right\rangle_q$$

Clearly \mathcal{F} is a lower bound of $\log p(y)$ and minimizing the KL divergence is equivalent to maximizing \mathcal{F} . Since $p(f, y) = p(y|f)p(f)$ we can equivalently write

$$\mathcal{F}(q(f)) = \langle \log p(y|f) \rangle_q + \langle \log p(f) \rangle_q - \langle \log q(f) \rangle_q$$

Gaussian Variational Approximation (GVA)

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

If we restrict $q \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$ and $p(f) \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$ a zero-mean Gaussian process we get

$$\begin{aligned}\mathcal{F}(q(f)) = & \langle \log p(y|f) \rangle_q + \frac{1}{2} \log |\mathbf{V}\mathbf{K}^{-1}| \\ & + \frac{n}{2} - \frac{1}{2} \mathbf{m}^T \mathbf{K}^{-1} \mathbf{m} - \frac{1}{2} \text{Tr}(\mathbf{V}\mathbf{K}^{-1})\end{aligned}$$

If the likelihood factorizes $p(y|f) = \prod p(y_i|f_i)$ then

$$\mathbf{V} = (\mathbf{K}^{-1} + \mathbf{\Lambda})^{-1}$$

with $\mathbf{\Lambda}$ diagonal (i.e., $2n$ variational parameters).

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- **GPML**: Gaussian Processes for Machine Learning
Rasmussen, C.E. and Nickisch, H. Gaussian
Processes for Machine Learning (GPML) Toolbox,
JMLR, 11 (2010) 3011-3015.
- **GPstuff**: Vanhatalo, J. Riihimäki, J. Hartikainen, J.
Jylänki, P. Tolvanen, V. and Vehtari, A. GPstuff:
Bayesian Modeling with Gaussian Processes. *JMLR*,
14 (2013) 1175-1179.

Tools

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

	GPstuff	GPML	FBM
Covariance functions			
number of elementary functions	13	10	4
sums of elements, masking of inputs	x	x	x
delta distance	x		x
products, positive scaling of elements	x	x	
Mean functions			
number of elementary functions	4	4	0
sums of elements, masking of inputs	x	x	
products, power, scaling of elements		x	
marginalized parameters	x		
Single latent likelihood/observation models			
Gaussian	x	x	x
logistic/logit, erf/probit	x	x	MCMC
Poisson	x	LA/EP/MCMC	MCMC
Gaussian scale mixture	MCMC		MCMC
Student- <i>t</i>	x	LA/VB/MCMC	
Laplacian		EP/VB/MCMC	
mixture of likelihoods		LA/EP/MCMC	
sech-squared, uniform for classification		x	
derivative observations	for sevp covf only		
binomial, negative binomial, zero-trunc. negative binomial, log-Gaussian	x		
Cox process; Weibull, log-Gaussian and log-logistic with censoring			
quantile regression	MCMC/EP		

Tools

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Multitask likelihood/observation models

multinomial, Cox proportional hazard model, density estimation, density regression, input dependent noise, input dependent overdispersion in Weibull, zero-inflated negative binomial	MCMC/LA		
multinomial logit (softmax)	MCMC/LA		MCMC
multinomial probit	EP		MCMC
Priors for parameters (ϑ)			
several priors, hierarchical priors	x		x
Sparse models			
FITC	x	exact/EP/LA	
CS, FIC, CS+FIC, PIC, VAR, DTC, SOR	x		
PASS-GP	LA/EP		
Latent inference			
exact (Gaussian only)	x	x	x
scaled Metropolis, HMC	x		x
LA, EP, elliptical slice sampling	x	x	
variational Bayes (VB)		x	
scaled HMC (with inverse of prior cov.)		x	
scaled HMC (whitening with approximate posterior covariance)	x		
parallel EP, Robust-EP	x		
marginal corrections (cm2 and fact)	x		

Tools

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Hyperparameter inference

type II ML	x	x	x
type II MAP, Metropolis, HMC	x		x
LOO-CV for Gaussian	x	x	
least squares LOO-CV for non-Gaussian		some likelihoods	
LA/EP LOO-CV for non-Gaussian, k-fold CV	x		
NUTS, slice sampling (SLS), surrogate SLS, shrinking-rank SLS, covariance-matching SLS, grid, CCD, importance sampling	x		

Model assessment

marginal likelihood	MAP, ML	ML
LOO-CV for fixed hyperparameters	x	x
LOO-CV for integrated hyperparameters, k-fold CV, WAIC, DIC	x	
average predictive comparison	x	

Transformation of Parameters

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

A constrained optimization problem, e.g., $\sigma^2 > 0$ can often be transformed into an unconstrained optimization problem. If $w = f(\theta)$ then

$$p_w(w) = |J|p_\theta(f^{-1}(w))$$

with J the Jacobian of the transformation between parameters

For example, if $w = \log(\sigma^2)$ then $w \in (-\infty, \infty)$ and

$$p_w(w) = |J|p_{\sigma^2}(\exp(w)) = \sigma^2 p_{\sigma^2}(\sigma^2)$$

Gaussian Process Applications

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

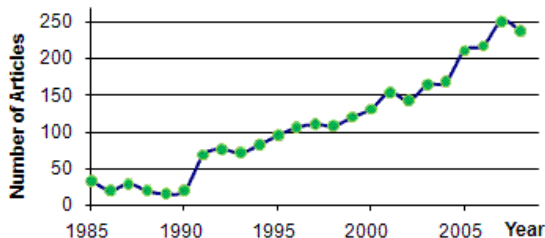
Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Clustering
- Ordinal Regression
- Preference Learning
- Ranking
- Surrogate modeling
- Global Optimization

- Relational Learning
- Reinforcement Learning
- Visualization high dim. data
- Nonrigid Shape Recovery
- Evaluating Integrals
- Dynamic systems



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Multiple Annotators

Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

- Instead of 1 annotator, we have M annotators.
- Items can be annotated by 1 or more annotators.
- Each annotator has its own noise level, expertise etc.
- How to combine annotations into a prediction?
- Recent growing interest in this type of problem (e.g., Mechanical Turk)

Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

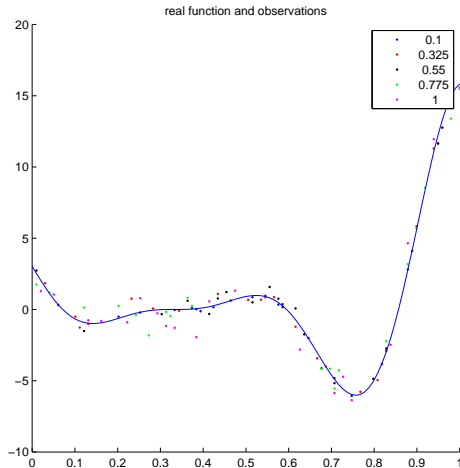
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



Multiple annotators

Let $X = [X_1, \dots, X_M]$, $Y = [Y_1, \dots, Y_M]$. A priori we can write $[Y, f_*] = [Y_1, \dots, Y_M, f_*]^T \sim$

$$\mathcal{N} \left(\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K(X_1, X_1) + \sigma_1^2 \mathbf{I}_1 & \cdots & K(X_1, X_M) \\ \vdots & & \vdots \\ K(X_M, X_1) & \cdots & K(X_M, X_M) + \sigma_M^2 \mathbf{I}_M \\ K(X_*, X_1) & \cdots & K(X_*, X_M) \end{bmatrix} \right)$$

with \mathbf{I}_m the $N_m \times N_m$ identity matrix. Let the diagonal matrix $\mathbf{N} = \text{diag}(\text{diag}(\sigma_1^2 \mathbf{I}_1), \dots, \text{diag}(\sigma_M^2 \mathbf{I}_M))$. The predictive equations are thus given by

$$\bar{f}_* = K(X_*, X) [K(X, X) + \mathbf{N}]^{-1} Y$$

$$\text{cov}(f_*) = K(X_*, X_*) - K(X_*, X) [K(X, X) + \mathbf{N}]^{-1} K(X, X_*)$$

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Let $X = \cup_{i=1}^M X_m$, $Y = [Y_1, \dots, Y_M]$. Define

$$\frac{1}{\hat{\sigma}_i^2} = \sum_{m \sim i} \frac{1}{\sigma_m^2},$$

$$\hat{y}_i = \hat{\sigma}_i^2 \sum_{m \sim i} \frac{y_i^m}{\sigma_m^2},$$

$$\hat{\Sigma} = \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_I^2)$$

with I the number of elements in X and $m \sim i$ denoting all annotators m that annotated x_i .

Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

The predictive equations are then given by

$$\bar{f}_* = K(X_*, X) \left[K(X, X) + \hat{\Sigma} \right]^{-1} \hat{Y}$$

$$\text{cov}(f_*) = K(X_*, X_*) - K(X_*, X) \left[K(X, X) + \hat{\Sigma} \right]^{-1} K(X, X_*)$$

and negative log-likelihood $-\log(Y)$

$$\begin{aligned} & \frac{1}{2} \log |K + \hat{\Sigma}| + \frac{1}{2} \hat{Y} (K + \hat{\Sigma})^{-1} \hat{Y} + \frac{N}{2} \log(2\pi) \\ & - \frac{1}{2} \log |\hat{\Sigma}| - \sum_i \sum_{m \sim i} \log \frac{1}{\sigma_m} + \frac{1}{2} \sum_i \sum_{m \sim i} \frac{(y_i^m)^2}{\sigma_m^2} - \frac{1}{2} \sum_i \frac{\hat{y}_i^2}{\hat{\sigma}_i^2} \end{aligned}$$

Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

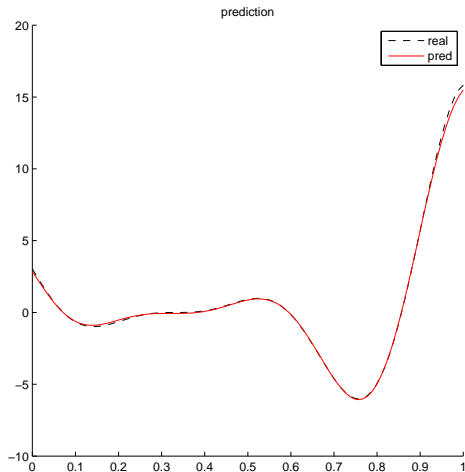
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

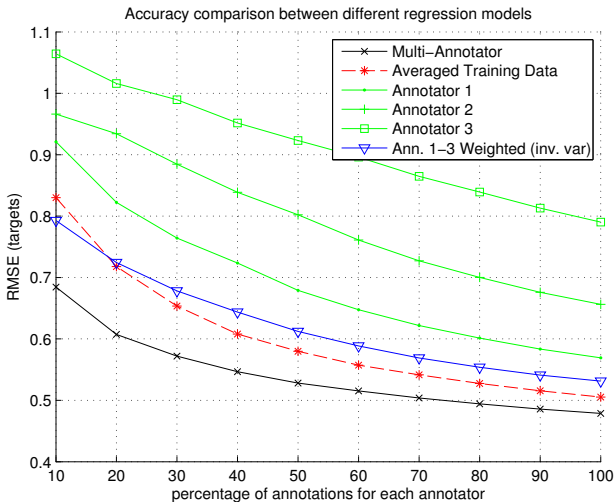
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



Multiple annotators

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

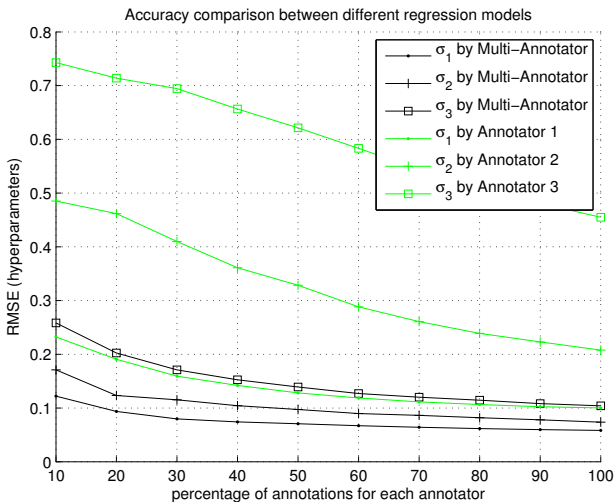
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Censored Observations

Censoring

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Goal is to learn a function

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

given a set of censored observations

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$$

where y is a censored version of y^* :

$$y = \begin{cases} l & \text{if } y^* \leq l \\ y^* & \text{if } l < y^* < u \\ u & \text{if } y^* \geq u \end{cases}$$

Censoring

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

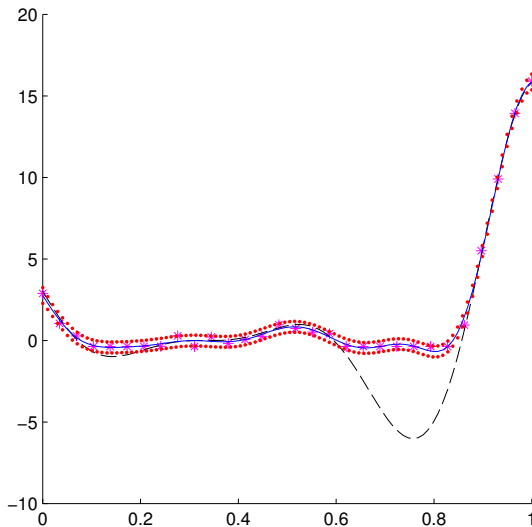
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Censoring

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

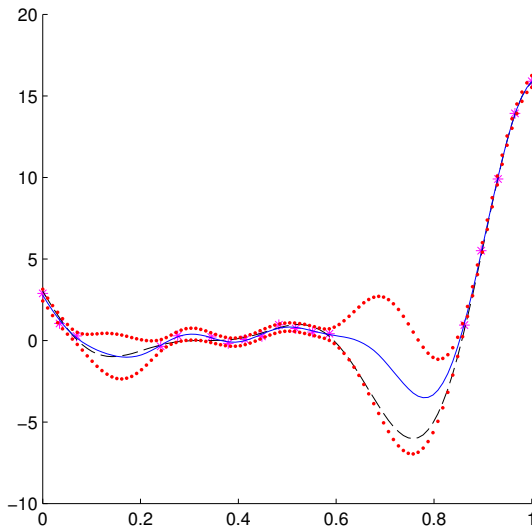
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Censoring

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Assume that latent function values are contaminated with Gaussian noise with zero mean and unknown variance.

Likelihood becomes a mixture of Gaussian and probit likelihood terms:

$$L = \prod_{i=1}^n p(y_i | f_i) = \prod_{y_i=l} \left[1 - \Phi \left(\frac{f_i - l}{\sigma} \right) \right] \prod_{l < y_i < u} \left[\frac{1}{\sigma} \phi \left(\frac{y_i - f_i}{\sigma} \right) \right] \prod_{y_i=u} \left[\Phi \left(\frac{f_i - u}{\sigma} \right) \right]$$

which is well-known as the **Tobit likelihood**.

Censoring

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

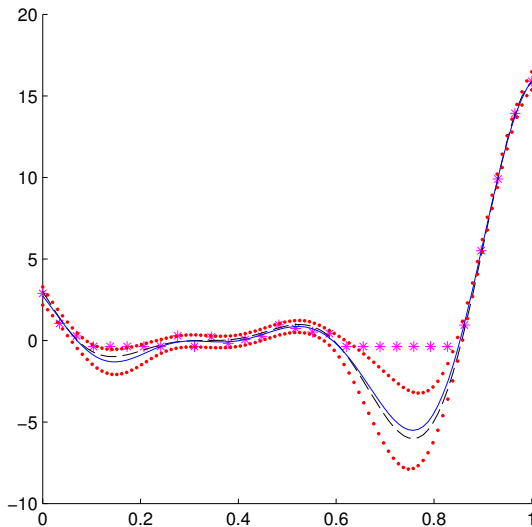
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Multi-step Prediction in Dynamic Systems

Dynamic Systems

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Consider a dynamical system

$$x_{t+1} = f(u_t, x_t) \quad y_t = x_t + \epsilon_t$$

with system state y and control u at time step t , and ϵ typically Gaussian white noise.

Dynamic Systems

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Training data can be obtained from system observations:

$$\Phi_N = \begin{bmatrix} u(1) & y(1) \\ u(2) & y(2) \\ \vdots & \vdots \\ u(k) & y(k) \\ \vdots & \vdots \\ u(N) & y(N) \end{bmatrix} \quad y_N = \begin{bmatrix} y(2) \\ y(3) \\ \vdots \\ y(k+1) \\ \vdots \\ y(N+1) \end{bmatrix}$$

Test data is constructed in similar fashion.

Naive multi-step prediction

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

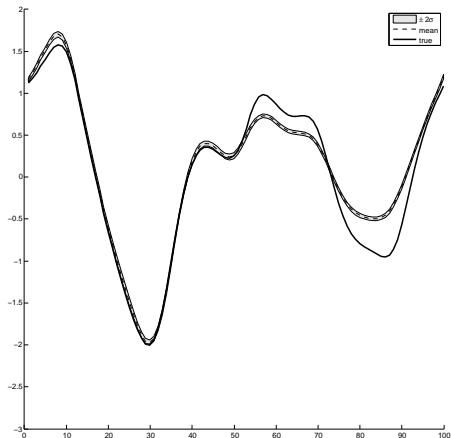
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Forecasting with Noisy Inputs

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Standard Gaussian process predictive equations:

$$\begin{aligned}\mu_* &= \mathbf{K}_{*N}(\mathbf{K}_{NN} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \\ \sigma_*^2 &= K_{**} - \mathbf{K}_{*N}(\mathbf{K}_{NN} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{N*} + \sigma^2\end{aligned}$$

If the input is given by a distribution $p(\mathbf{x}_*|m, V) \sim \mathcal{N}(m, V)$, the predictive distribution is given by integrating over the input distribution:

$$p(f_*|m, V, D) = \int p(f_*|\mathbf{x}_*, D)p(\mathbf{x}_*|m, V) d\mathbf{x}_*$$

Forecasting with Uncertainty Propagation

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

$$p(f_* | m_t, V_t, D) = \underbrace{\int \underbrace{p(f_* | x_*, D)}_{\text{nonlinear}} \underbrace{p(x_* | m_t, V_t)}_{\text{Gaussian}} dx_*}_{\text{non-Gaussian}}$$

$$\approx \mathcal{N}(m_{t+1}, V_{t+1})$$

$$p(f_* | m_{t+1}, V_{t+1}, D) = \dots$$

Propagating Uncertainty

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Input distribution is given by $p(x_{T+n} | Y_T) \sim \mathcal{N}(u_{T+n}, S_{T+n})$.

At $t = T + 1$

$$u_{T+1} = [y_{T+1-L}, \dots, y_T],$$

$$S_{T+1} = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix}$$

and

$$p(y_{T+1} | Y_T) \sim \mathcal{N}(\mu(u_{T+1}), \sigma^2(u_{T+1}) + \sigma_\epsilon^2)$$

Propagating Uncertainty

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

At $t = T + 2$

$$u_{T+2} = [y_{T+2-L}, \dots, y_T, \mu(u_{T+1})]$$

$$S_{T+2} = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \sigma^2(u_{T+1}) + \sigma_\epsilon^2 \end{bmatrix}$$

and

$$p(y_{T+2} | Y_T) \sim \mathcal{N}(m(u_{T+2}, S_{T+2}), v(u_{T+2}, S_{T+2}) + \sigma_\epsilon^2)$$

Propagating Uncertainty

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

At $t = T + k$

$$u_{T+k} = [m(u_{T+k-L}, S_{T+k-L}), \dots, m(u_{T+k-1}, S_{T+k-1})]$$

$$S_{T+k} = \begin{bmatrix} v(u_{T+k-L}, S_{T+k-L}) + \sigma_{\epsilon}^2 & \dots & \text{cov}(y_{T+k-L}, y_{T+k-1}) \\ \vdots & & \vdots \\ \text{cov}(y_{T+k-L}, y_{T+k-1}) & \dots & v(u_{T+k-1}, S_{T+k-1}) + \sigma_{\epsilon}^2 \end{bmatrix}$$

A GP with Gaussian kernel and Gaussian input distribution allows m_{t+1} , V_{t+1} to be computed *analytically*.

Dynamic Systems

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

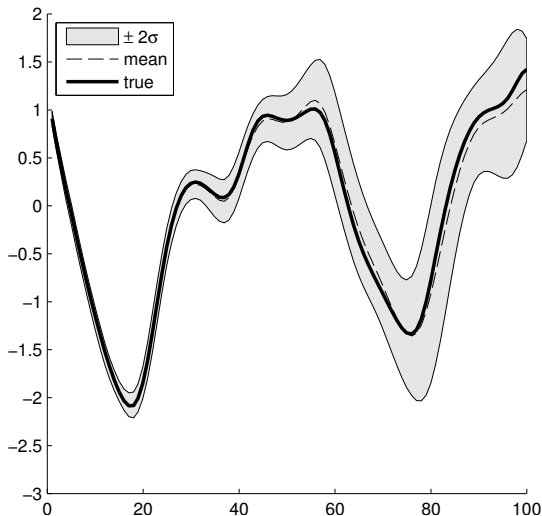
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Dynamic Systems

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

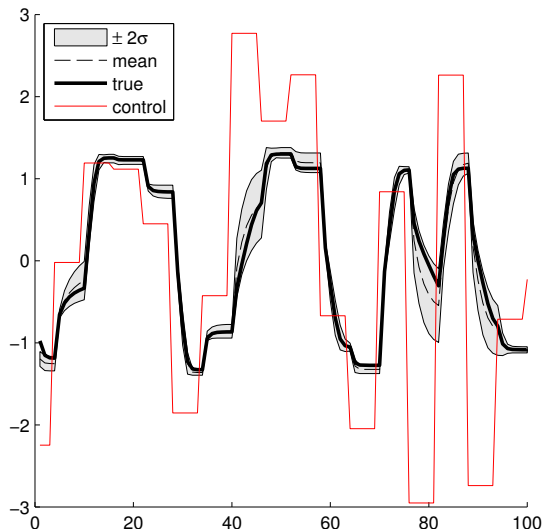
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

High- and Low Fidelity Observations

Multi-fidelity Analysis

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

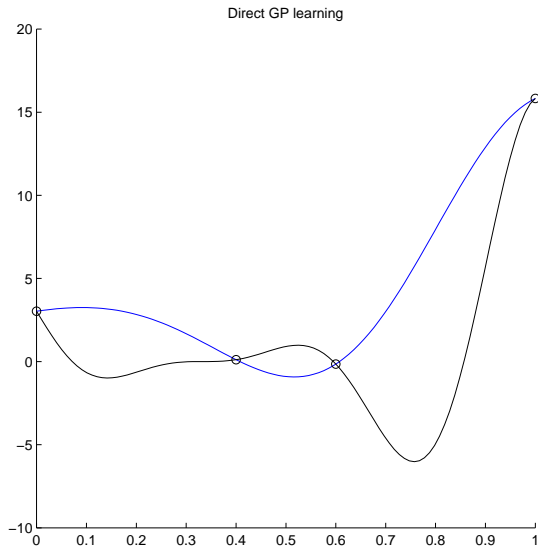
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Multi-fidelity Analysis

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

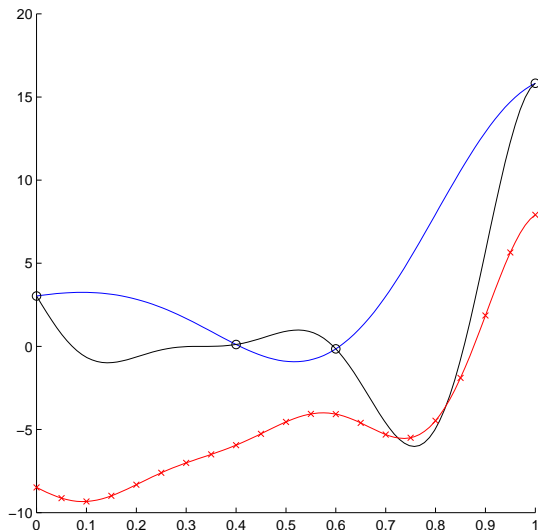
Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning



Multi-fidelity Analysis

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

We will use the following model:

$$y_h(x) = \rho y_l(x) + d(x)$$

with ρ a **scaling** parameter and $d(x)$ a GP modeling the **difference** $y_h(x) - \rho y_l(x)$. We assume

$$\text{cov}\{Y_h(x'), Y_l(x) | Y_l(x')\} = 0, \forall x \neq x'$$

and l, d independent GPs.

Multi-fidelity Analysis

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

The covariance for the high-fidelity model can then be computed on $Y = [Y_l(X_l), Y_h(X_h)]$ using

$$\begin{aligned}\text{cov}\{Y_l(X_l), Y_l(X_l)\} &= \text{cov}\{I(X_l), I(X_l)\} = K_l(X_l, X_l) \\ \text{cov}\{Y_h(X_h), Y_l(X_l)\} &= \text{cov}\{\rho I(X_h) + d(X_h), I(X_l)\} \\ &= \rho \text{cov}\{I(X_h), I(X_l)\} = \rho K_l(X_h, X_l) \\ \text{cov}\{Y_h(X_h), Y_h(X_h)\} &= \text{cov}\{\rho I(X_h) + d(X_h), \rho I(X_h) + d(X_h)\} \\ &= \rho^2 \text{cov}\{I(X_h), I(X_h)\} + \text{cov}\{d(X_h), d(X_h)\} \\ &= \rho^2 K_l(X_h, X_h) + K_d(X_h, X_h)\end{aligned}$$

Multi-fidelity Analysis

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

The covariance matrix K_h is thus given by

$$K_h = \begin{bmatrix} K_l(X_l, X_l) & \rho K_l(X_l, X_h) \\ \rho K_l(X_h, X_l) & \rho^2 K_l(X_h, X_h) + K_d(X_h, X_h) \end{bmatrix}$$

and we can make predictions

$$\begin{aligned} \overline{y_{h*}} &= K_h(X_*, X) K_h(X, X)^{-1} Y \\ \text{cov}(y_{h*}) &= K_h(X_*, X_*) - K_h(X_*, X) K_h(X, X)^{-1} K_h(X, X_*) \end{aligned}$$

Multi-fidelity Analysis

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

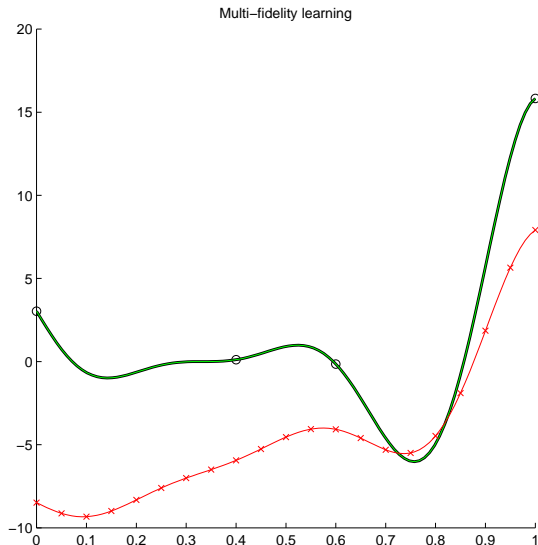
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Gaussian Process Surrogate Models

Surrogate Modelling with GPs

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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

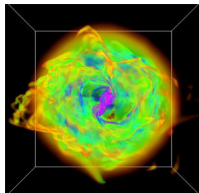
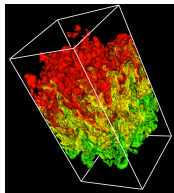
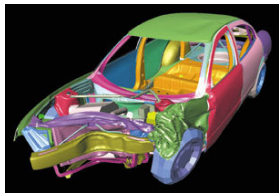


Figure: 1. Car collision; 2. Turbulent-mixing dynamics of a supernova; 3. Gas cloud collapsing inwards to form a star.

Surrogate Modelling with GPs

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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).

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

Function Optimization

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

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

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

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

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

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

EI - 1D example

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

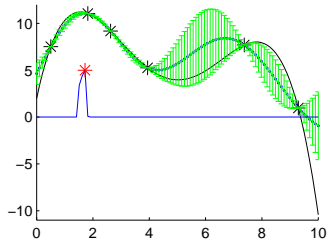
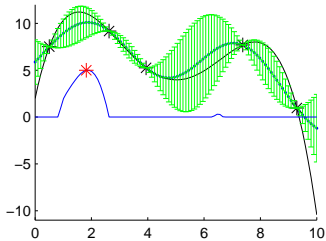
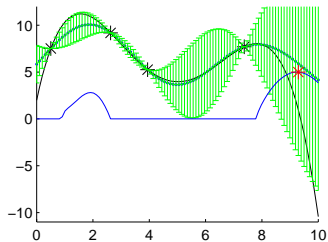
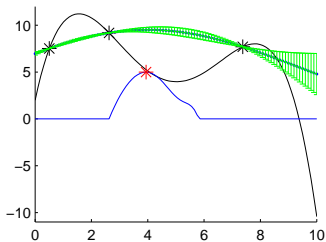
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Gaussian Process Integral Predictions

Bayesian Monte Carlo and Optimization

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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).



Problem Setting

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

Assume some underlying, unknown real-valued function, that can be evaluated

$$f(\mathbf{x}_c, \mathbf{x}_e) \rightarrow \mathbb{R}$$

Our optimization problem can be formulated as

$$\begin{aligned}\mathbf{x}_c^* &= \operatorname{argmax}_{\mathbf{x}_c} E[\ell(\mathbf{x}_c)] \\ &= \operatorname{argmax}_{\mathbf{x}_c} \int_{\mathbf{x}_e} f(\mathbf{x}_c, \mathbf{x}_e) p(\mathbf{x}_e) d\mathbf{x}_e\end{aligned}$$

Bayesian Monte Carlo

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

We can think of F as being **random** as we are uncertain about $f(\mathbf{x})$ because we have a limited number of samples [O'Hagan, 1991; Rasmussen & Ghahramani, 2003].

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)

Bayesian Monte Carlo

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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.

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

$$\begin{aligned} p(\mathbf{x}) &\sim \mathcal{N}(\mathbf{b}, \mathbf{B}) \\ k(\mathbf{x}, \mathbf{x}') &= w_0 \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \mathbf{A}^{-1}(\mathbf{x} - \mathbf{x}')\right) \end{aligned}$$

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

Bayesian Monte Carlo - 1D demo

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

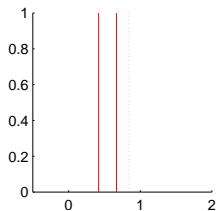
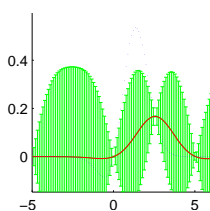
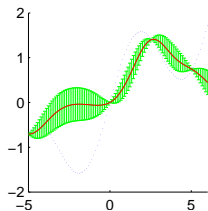
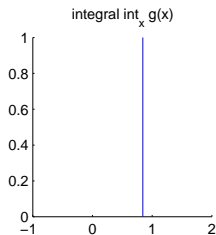
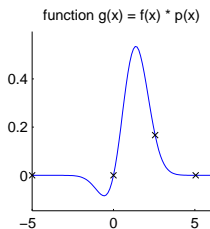
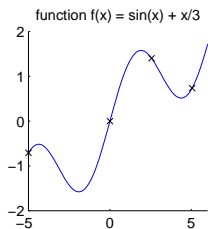
Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning



Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

Preference Learning

Preference Learning

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

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)

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

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 $\mathbf{x}_1 \succ \mathbf{x}_2$.

Preference Learning

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior

Sampling

Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators

Censoring

Dynamic Systems

Multi-fidelity Analysis

Surrogate Modeling

Integration

Preference Learning

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

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

Preference Learning

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

Perry Groot

Regression

Probabilistic
Inference

Gaussian
processes

Posterior
Sampling
Model Selection

Classification

Approximations

Tools

Applications

Multiple annotators
Censoring
Dynamic Systems
Multi-fidelity Analysis
Surrogate Modeling
Integration
Preference Learning

