

GAUSSIAN PROCESSES WITH APPLICATIONS

MATTIAS VILLANI

**DEPARTMENT OF STATISTICS
STOCKHOLM UNIVERSITY**

AND

**DEPARTMENT OF COMPUTER AND INFORMATION SCIENCE
LINKÖPING UNIVERSITY**

- Introduction to **Gaussian process regression**
- Gaussian process all the things - some **applications**
- No time for **Bayesian optimization** :(

COOL IN MACHINE LEARNING NOW, BUT ... STATISTICIANS WERE ROCKIN' IT ALREADY IN THE 70'S

Biometrika (1975), **62**, 1, p. 79

With 5 text-figures

Printed in Great Britain

79

A Bayesian approach to model inadequacy for polynomial regression

BY B. J. N. BLIGHT

Department of Statistics, Birkbeck College, London

AND L. OTT

Department of Statistics, University of Florida, Gainesville

COOL IN MACHINE LEARNING NOW, BUT ... STATISTICIANS WERE ROCKIN' IT ALREADY IN THE 70'S

J. R. Statist. Soc. B (1978),
40, No. 1, pp. 1-42

Curve Fitting and Optimal Design for Prediction

By A. O'HAGAN

University of Warwick

[Read before the ROYAL STATISTICAL SOCIETY at a meeting organized by the RESEARCH SECTION
on Wednesday, October 12th, 1977, Professor J. F. C. KINGMAN in the Chair]

SUMMARY

A Bayesian approach to density estimation

THORBURN DANIEL

Biometrika, Volume 73, Issue 1, 1 April 1986, Pages 65–75, <https://doi-org.e.bibl.liu.se/10.1093/biomet/73.1.65>

Published: 01 April 1986 **Article history** ▼

■ Linear regression

$$y = f(\mathbf{x}) + \epsilon$$

$$f(\mathbf{x}) = \mathbf{x}^T \beta$$

and $\epsilon \sim N(0, \sigma_n^2)$ and iid over observations.

■ Polynomial regression: $\phi(\mathbf{x}) = (1, x, x^2, x^3, \dots, x^k)$:

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \beta.$$

■ More generally: **splines** with **basis functions**.

■ Polynomial and spline models are linear in β . Least squares!

- **Model:** Linear regression for all n observations

$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\beta} + \underset{n \times 1}{\varepsilon}$$

- **Prior**

$$\beta \sim N(0, \Sigma_p)$$

- Common choice (Ridge regression): $\Sigma_p = \lambda^{-1} \mathbf{I}$.

- **Posterior**

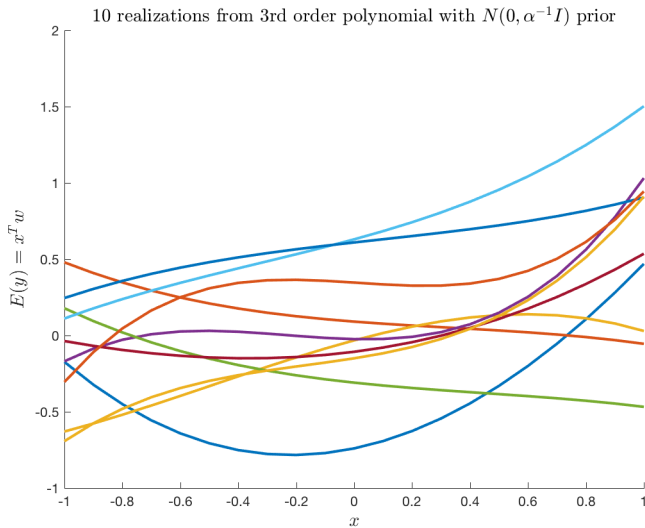
$$\beta | \mathbf{X}, \mathbf{y} \sim N(\bar{\beta}, \mathbf{A}^{-1})$$

$$\mathbf{A} = \sigma_n^{-2} \mathbf{X}^T \mathbf{X} + \Sigma_p^{-1}$$

$$\bar{\beta} = \sigma_n^{-2} \left(\sigma_n^{-2} \mathbf{X}^T \mathbf{X} + \Sigma_p^{-1} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

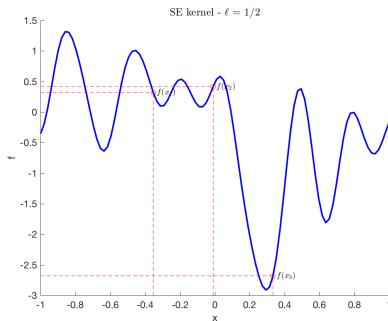
- **Posterior precision = Data Precision + Prior Precision.**

A PRIOR ON β IS REALLY A PRIOR OVER FUNCTIONS



NON-PARAMETRIC REGRESSION

- **Non-parametric regression:** avoid a parametric form for $f(\cdot)$.
- Treat $f(\mathbf{x})$ as **an unknown parameter for every \mathbf{x}** .



- A *new* parameter for *every* \mathbf{x} , you must be joking?
- Instead of restricting to linear, impose **smoothness**.

■ Weight space view

- Restrict attention to a grid of x-values: x_1, \dots, x_k .
- Put a joint prior on the **vector of k function values**

$$f(x_1), \dots, f(x_k)$$

■ Function space view

- Treat **f as an unknown function.**
- Put a prior over a set of functions.

- A GP implies:

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_k) \end{pmatrix} \sim N(\mathbf{m}, \mathbf{K})$$

- But how do we specify the $k \times k$ **covariance matrix** \mathbf{K} ?

$$\text{Cov}(f(x_p), f(x_q))$$

- **Squared exponential covariance function**

$$\text{Cov}(f(x_p), f(x_q)) = k(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2} \left(\frac{x_p - x_q}{\ell}\right)^2\right)$$

- Nearby x 's have highly correlated function ordinates $f(x)$.
- We can compute $\text{Cov}(f(x_p), f(x_q))$ for *any* x_p and x_q .

Definition

A **Gaussian process (GP)** is a collection of random variables, any finite number of which have a multivariate Gaussian distribution.

- A GP is a **probability distribution over functions**.
- A GP is specified by a **mean** and a **covariance function**

$$m(x) = E[f(x)]$$

$$k(x, x') = E[(f(x) - m(x))(f(x') - m(x'))]$$

for any two inputs x and x' .

- A **Gaussian process** is denoted by

$$f(x) \sim \text{GP}(m(x), k(x, x'))$$

- $f(x) \sim \text{GP}$ encodes **prior beliefs** about the unknown $f(\cdot)$.

■ Let $r = \|x - x'\|$.

■ **Squared exponential (SE)** kernel ($\ell > 0, \sigma_f > 0$)

$$K_{SE}(r) = \sigma_f^2 \exp\left(-\frac{r^2}{2\ell^2}\right)$$

■ **Matérn** kernel ($\ell > 0, \sigma_f > 0, \nu > 0$)

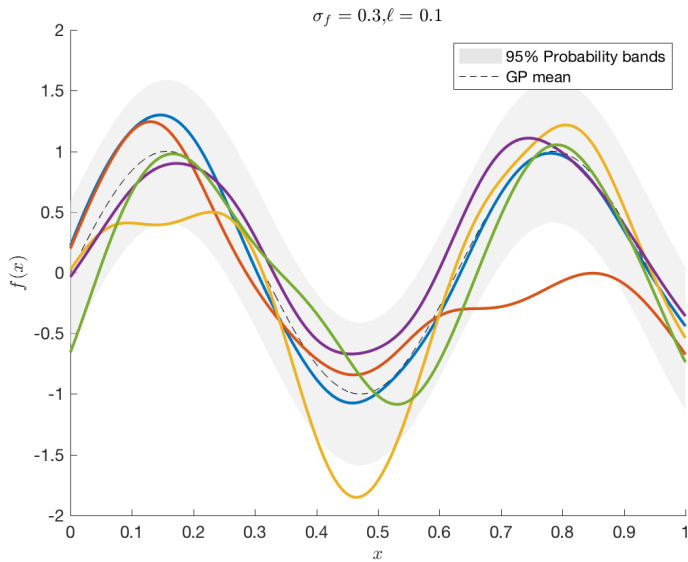
$$K_{Matern}(r) = \sigma_f^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\ell}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu}r}{\ell}\right)$$

■ **Simulate draw** from $f(x) \sim \text{GP}(m(x), k(x, x'))$ by:

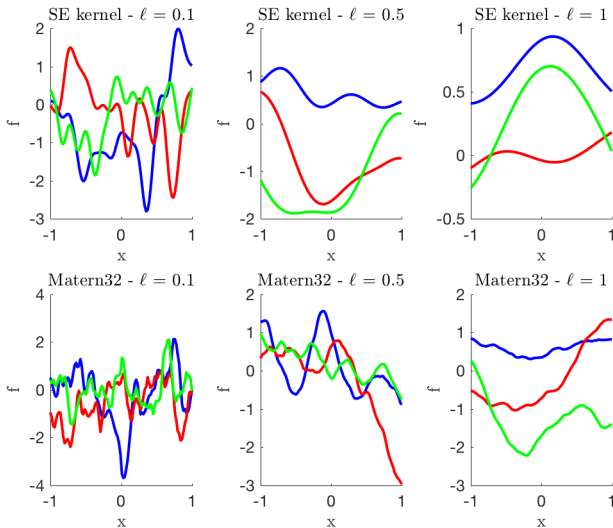
- form a grid $\mathbf{x}_* = (x_1, \dots, x_n)$
- simulate function values from multivariate normal:

$$f(\mathbf{x}_*) \sim N(m(\mathbf{x}_*), K(\mathbf{x}_*, \mathbf{x}_*))$$

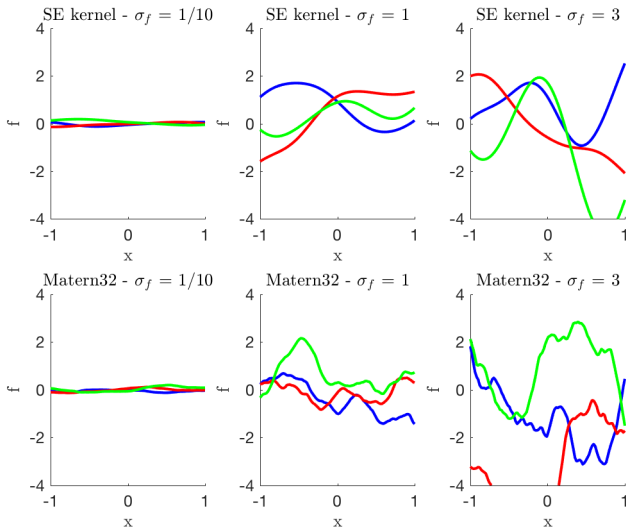
SIMULATING A GP



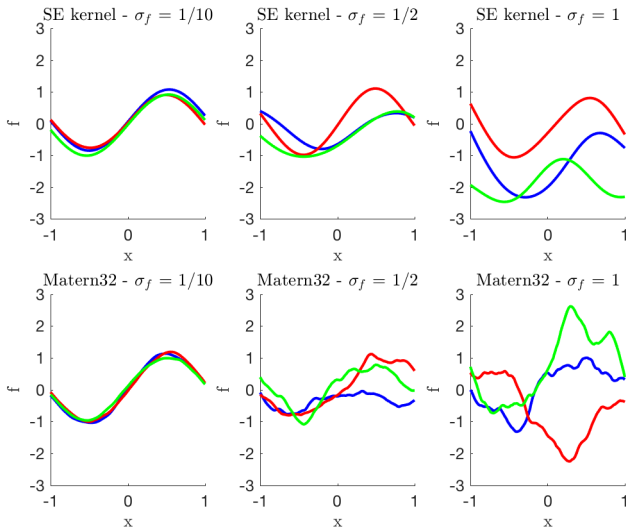
THE LENGTH SCALE ℓ DETERMINES THE SMOOTHNESS



THE SCALE FACTOR σ_f DETERMINES THE VARIANCE



THE MEAN CAN BE $\sin(3x)$. OR WHATEVER.



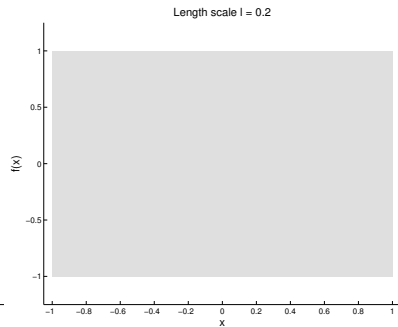
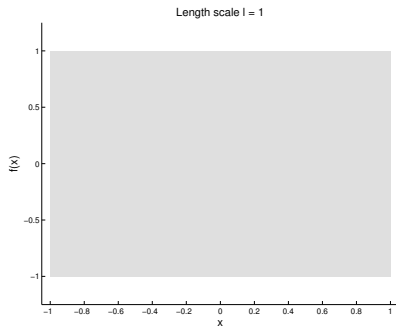
- The joint way: Choose a grid x_1, \dots, x_k . Simulate the k -vector

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_k) \end{pmatrix} \sim N(\mathbf{m}, \mathbf{K})$$

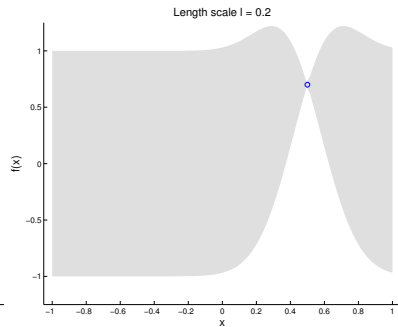
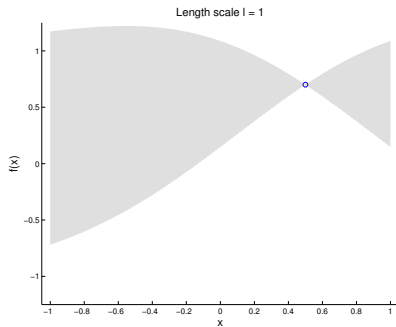
- More intuition from the conditional decomposition

$$\begin{aligned} p(f(x_1), f(x_2), \dots, f(x_k)) &= p(f(x_1)) p(f(x_2)|f(x_1)) \cdots \\ &\quad \times p(f(x_k)|f(x_1), \dots, f(x_{k-1})) \end{aligned}$$

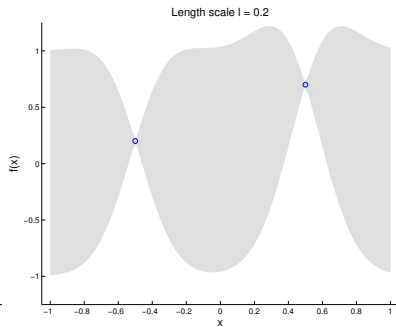
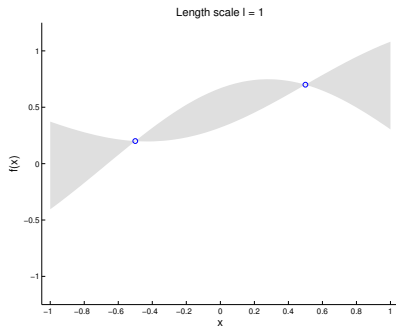
SIMULATING FROM $p(f(x_1))$



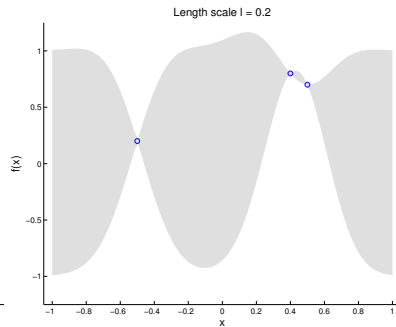
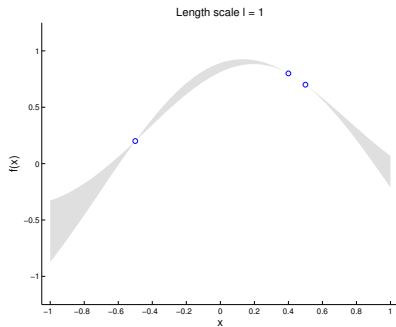
SIMULATING FROM $p(f(x_2)|f(x_1))$



SIMULATING FROM $p(f(x_3)|f(x_1), f(x_2))$



SIMULATING FROM $p(f(x_4)|f(x_1),f(x_2),f(x_3))$



THE POSTERIOR FOR A GAUSSIAN PROCESS REGRESSION

■ Model

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma_n^2)$$

■ Prior

$$f(x) \sim GP(0, k(x, x'))$$

■ **Observed:** $\mathbf{x} = (x_1, \dots, x_n)^T$ and $\mathbf{y} = (y_1, \dots, y_n)^T$.

■ **Goal:** posterior of $f(\cdot)$ over a grid of x -values: $\mathbf{f}_* = \mathbf{f}(\mathbf{x}_*)$.

■ Posterior

$$\mathbf{f}_* | \mathbf{x}, \mathbf{y}, \mathbf{x}_* \sim N(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*))$$

$$\bar{\mathbf{f}}_* = K(\mathbf{x}_*, \mathbf{x}) [K(\mathbf{x}, \mathbf{x}) + \sigma_n^2 I]^{-1} \mathbf{y}$$

$$\text{cov}(\mathbf{f}_*) = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, \mathbf{x}) [K(\mathbf{x}, \mathbf{x}) + \sigma_n^2 I]^{-1} K(\mathbf{x}, \mathbf{x}_*)$$

- Idea: obtain joint $p(\mathbf{y}, \mathbf{f}_*)$ and then $p(\mathbf{f}_*|\mathbf{y})$ by conditioning.

- **Model**

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma_n^2)$$

- **Prior**

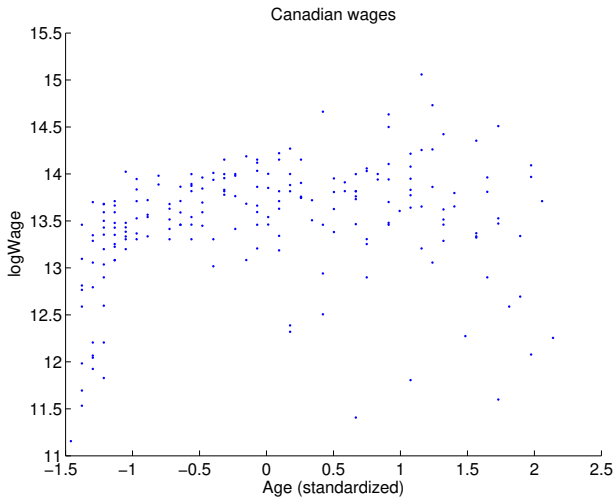
$$f(x) \sim GP(0, k(x, x'))$$

- Joint distribution of $(\mathbf{y}, \mathbf{f}_*)$

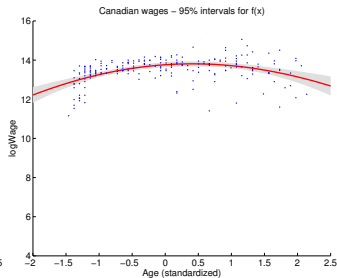
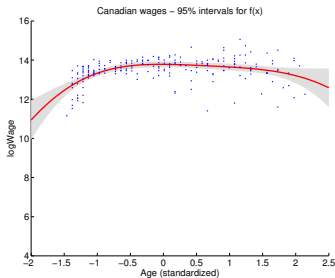
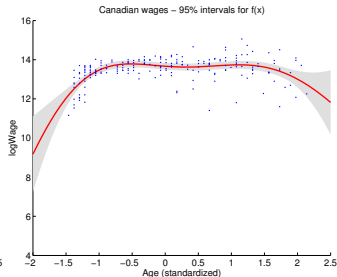
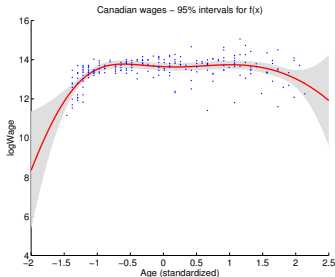
$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim N \left[\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) + \sigma_n^2 I & K(\mathbf{x}, \mathbf{x}_*) \\ K(\mathbf{x}_*, \mathbf{x}) & K(\mathbf{x}_*, \mathbf{x}_*) \end{pmatrix} \right]$$

- Result: conditional distributions from multivariate normal are normal.

EXAMPLE - CANADIAN WAGES



POSTERIOR OF F - $\ell = 0.2, 0.5, 1, 2$



- Kernel depends on **hyperparameters** θ . Example SE kernel $[\theta = (\sigma_f, \ell)^T]$

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left(-\frac{1}{2} \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\ell^2} \right)$$

- Common: maximize the **marginal likelihood** wrt θ :

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

$\mathbf{f} = f(\mathbf{X})$ is a vector of function values in the training data.

- For **Gaussian process regression**:

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

- Proper **Bayesian inference for hyperparameters**

$$p(\theta|\mathbf{y}, \mathbf{X}) \propto p(\mathbf{y}|\mathbf{X}, \theta) p(\theta).$$

- **Binary** or multi-class **response**. Aim: $\Pr(y_i = 1|\mathbf{x}_i)$.
- **Logistic regression**

$$\Pr(y_i = 1|\mathbf{x}_i) = \lambda(\mathbf{x}_i^T \beta), \text{ where } \lambda(z) = \frac{1}{1 + \exp(-z)}.$$

- $\lambda(z)$ 'squashes' the linear prediction $\mathbf{x}^T \beta \in \mathbb{R}$ into $[0, 1]$.
- **Linear decision boundaries** because of linear predictor $\mathbf{x}^T \beta$.
- **GP classification**: replace $\mathbf{x}^T \beta$ by $f(\mathbf{x})$ where

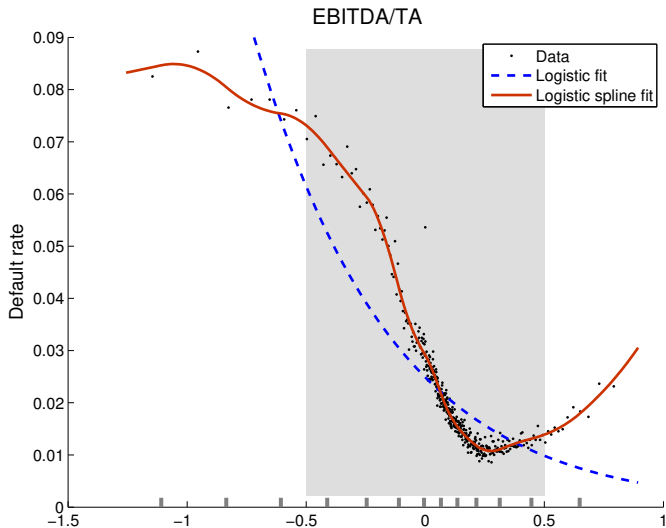
$$f \sim \text{GP}(0, k(\mathbf{x}, \mathbf{x}'))$$

and squash f through logistic function

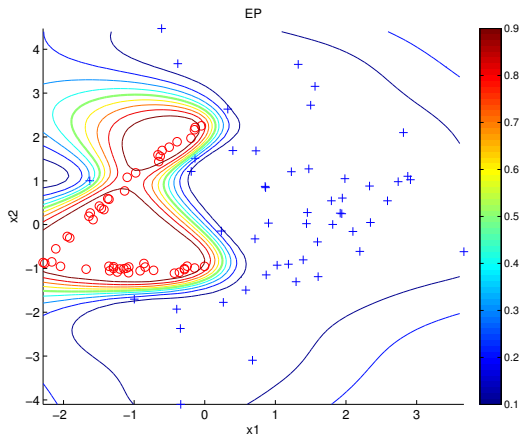
$$\Pr(y = 1|\mathbf{x}) = \lambda(f(\mathbf{x}))$$

- Nonparametric **flexible decision boundaries**.

CLASSIFICATION - FIRM BANKRUPTCY



GP CLASSIFICATION ON SIMULATED DATA



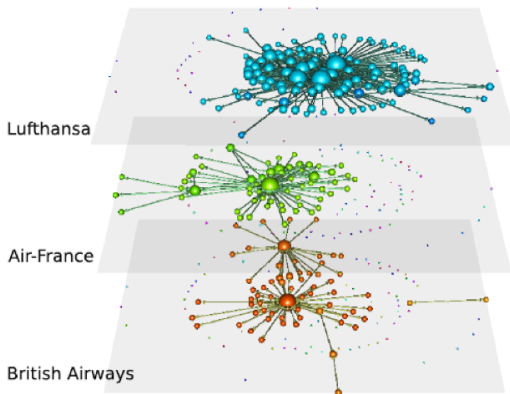
URBAN TRAFFIC PREDICTION

- **Aim:** Real-time prediction of road travel times.
- **Data:** noisy GPS from a large number of taxis. SL bus data.
- **Model:** **Structured Multivariate Gaussian Process:**
 - **GP**s over time for each street section. Periodic kernels.
 - **Hierarchical** factor-like model. Traffic on main roads feed the small roads.
 - **Spatial, topological**, dependence between street sections.



AIRLINE NETWORK PREDICTIONS

- **Aim:** **Predict** the evolution of airline **networks over time**.
- **Data:** Quarterly world-wide networks for all airlines.
- **Model:** **Dynamic multi-layered networks** driven by **GPs**.



- **Bernoulli model** with **bilinear latent Gaussian processes**:

$$Y_{uv}^{(k)}(t) | \pi_{uv}^{(k)}(t) \sim \text{Bern} \left[\pi_{uv}^{(k)}(t) \right]$$

$$\text{Logit} \left[\pi_{uv}^{(k)}(t) \right] = \mu(t) + \sum_{r=1}^R \bar{x}_{ur}(t) \bar{x}_{vr}(t) + \sum_{h=1}^H x_{uh}^{(k)}(t) x_{vh}^{(k)}(t),$$

where

- **Global GP** across actors and layers: $\mu(t)$
- **Actor-specific GPs**, but common across layers: $\bar{x}_{ur}(t)$.
- **Actor-specific and layer-specific** GPs: $x_{uh}^{(k)}(t)$.

- Scalability.

Thanks!