

Adventures in Gaussian Processes

AI4Research seminar, Uppsala University

Mattias Villani

Department of Statistics
Stockholm University

Department of Computer and Information Science
Linköping University

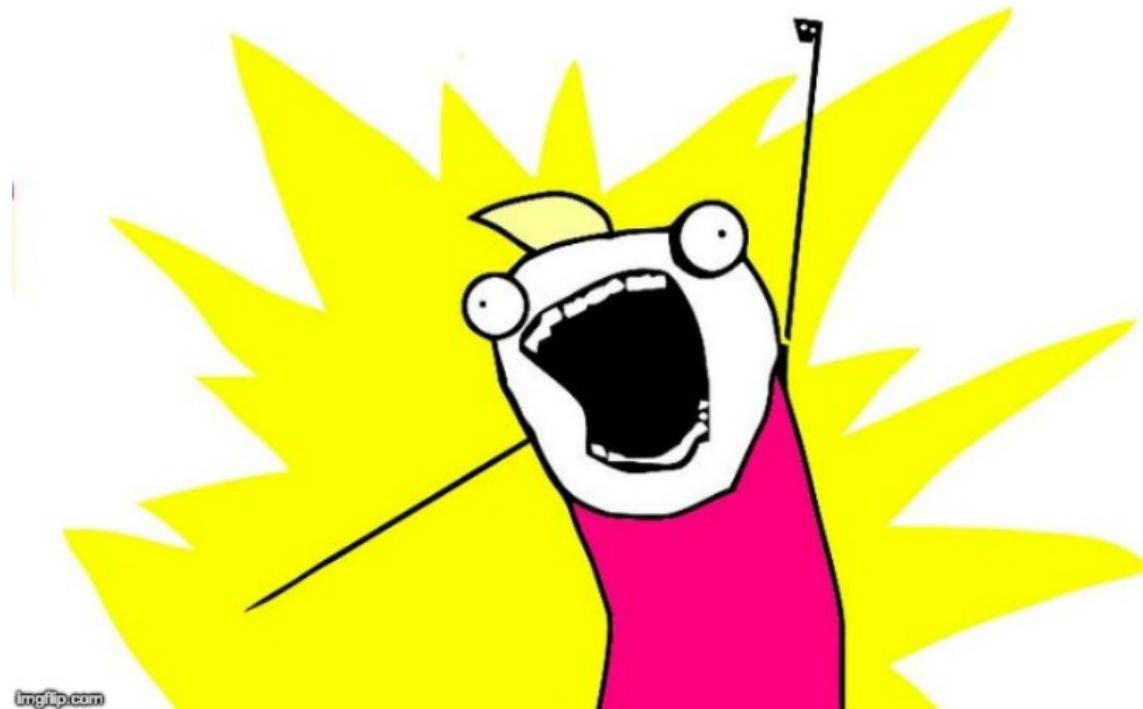


Overview

- Gaussian processes for data analysis - brief intro.
- Search-and-rescue using hierarchical spatial point processes.
- Dynamic multi-layer network modeling and prediction of flight connections.
- [Bayesian optimization with user-controlled precision.]

Alternative title

GAUSSIAN PROCESS ALL THE THINGS



Gaussian process regression

■ Gaussian process regression

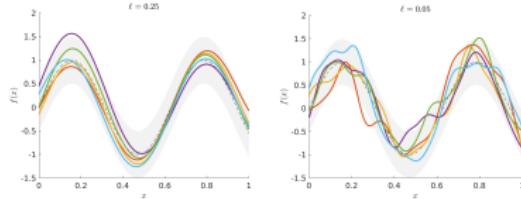
$$y_i = f(x_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma_n^2)$$

■ Gaussian process prior over the space of functions

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

with **squared exponential** covariance function

$$k(x, x') \equiv \text{Cov}(f(x), f(x')) = \sigma_f^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right)$$



Gaussian process regression

■ Model

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

■ Prior

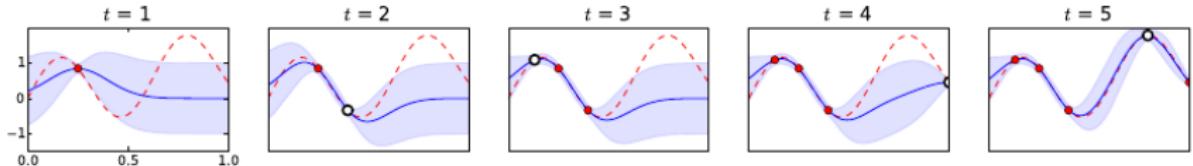
$$f(\mathbf{x}) \sim GP(0, k(\mathbf{x}, \mathbf{x}'))$$

■ Posterior of $f(\cdot)$ over test data \mathbf{X}_* :

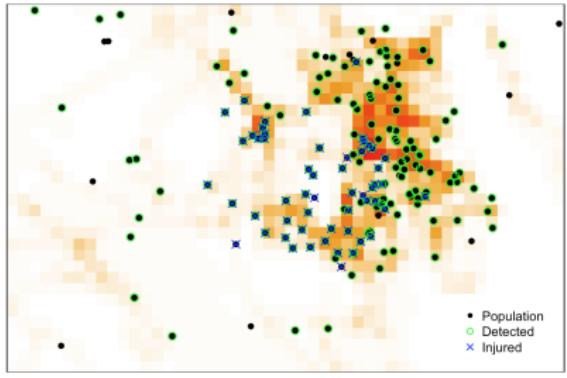
$$f_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_* \sim N(\bar{\mathbf{f}}_*, \text{cov}(\bar{\mathbf{f}}_*))$$

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

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



Search-and-rescue using spatial point processes¹



¹Andersson et al (2019). Real-Time Robotic Search using Structural Spatial Point Processes. *UAI*.

Hierarchical Spatial Point Process

- Log Gaussian Cox Process (LGCP) for number of persons in the region $\tilde{S} \subset S$.

$$N_{y^*}(\tilde{S}) | \lambda \sim \text{Poisson} \left(\int_{\mathbf{s} \in \tilde{S}} \lambda(\mathbf{s}) d\mathbf{s} \right)$$

$$\log \lambda(\mathbf{s}) = \alpha_\lambda + \underbrace{\mathbf{x}_\lambda^\top(\mathbf{s}) \boldsymbol{\beta}_\lambda}_{GIS} + \underbrace{\xi_\lambda(\mathbf{s})}_{GP \text{ in } 2D}$$

- The number of detected persons by a thinned LGCP

$$N_y(\tilde{S}) | r, \lambda \sim \text{Poisson} \left(\int_{\mathbf{s} \in \tilde{S}} r(\mathbf{s}) \lambda(\mathbf{s}) d\mathbf{s} \right)$$

$$\log r(\mathbf{s}) = \mathbf{x}_r^\top(\mathbf{s}) \boldsymbol{\beta}_r$$

- Probability of injury for i th person at location \mathbf{y}_i

$$w_i | q \sim \text{Bernoulli}(q(\mathbf{y}_i)),$$

$$\log \left(\frac{q(\mathbf{s})}{1 - q(\mathbf{s})} \right) = \alpha_q + \mathbf{x}_q^\top(\mathbf{s}) \boldsymbol{\beta}_q + \xi_q(\mathbf{s})$$

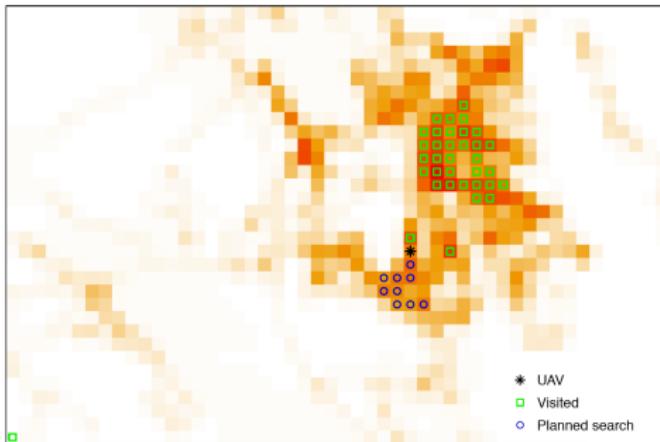
Real-time decision making under uncertainty

Challenges

- 1 **missing data** - point pattern is only partially observed
- 2 **real-time sequential inference over spatial fields**
- 3 **real-time decision making** under uncertainty

Solutions

- 1 Strong priors based on **GIS data**
- 2 **Warm-started INLA** for inference
- 3 Tailored **Monte Carlo Tree Search** for decision



Discretization of the process

- Partition spatial domain in $a_1 \cdot a_2$ rectangles with area Δ .
- **Number of detected** in region (i, j) :

$$n_{ij} \sim \text{Pois}(\Delta \exp(z_{\lambda,ij} + z_{r,ij}))$$

- Different features in population and detection models.
- **Number of detected injured** in region (i, j) :

$$m_{ij}|n_{ij} \sim \text{Binom}\left(n_{ij}, \frac{1}{1 + \exp(-z_{q,ij})}\right)$$

- The latent variables $z_{\lambda,ij}$, $z_{r,ij}$ and $z_{q,ij}$ are representative values of $\log \lambda(\mathbf{s})$, $\log r(\mathbf{s})$ and $\text{logit } q(\mathbf{s})$ in region s_{ij} .

Integrated Nested Laplace Approximation (INLA)

- INLA targets models of the form

$$\begin{aligned}\theta &\sim \pi(\theta) \\ z|\theta &\sim N(0, Q^{-1}(\theta)), \quad Q \text{ sparse}\end{aligned}$$

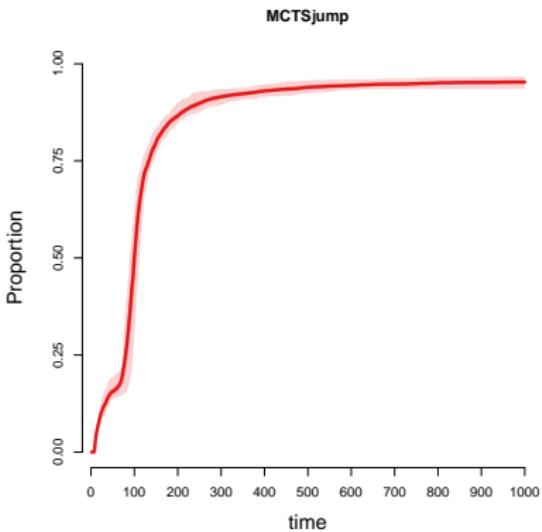
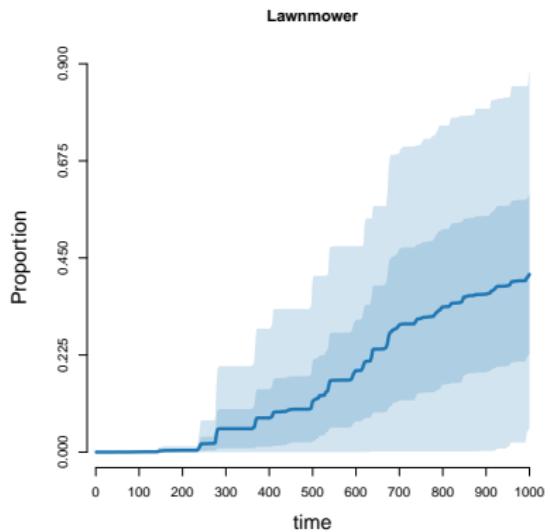
$$\pi(y|z, \theta) = \prod_{i=1}^n \pi(y_i|z_i, \theta)$$

- n is the number of rectangles
- y_i is number of detected or injured in rectangle i
- z_i is $z_{\lambda,ij} + z_{r,ij}$ or $z_{q,ij}$ in rectangle i (+ fixed effects α and β)
- θ is a low-dim vector with kernel hyperparameters
- INLA approximates marginal posterior of latents

$$\pi(z_i|y) = \int \tilde{\pi}(z_i|\theta, y) \tilde{\pi}(\theta|y) d\theta$$

by Laplace/Gaussian approx (z) and numerical integration (θ).

We find injured a lot faster than lawnmower



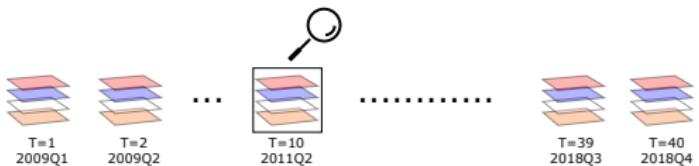
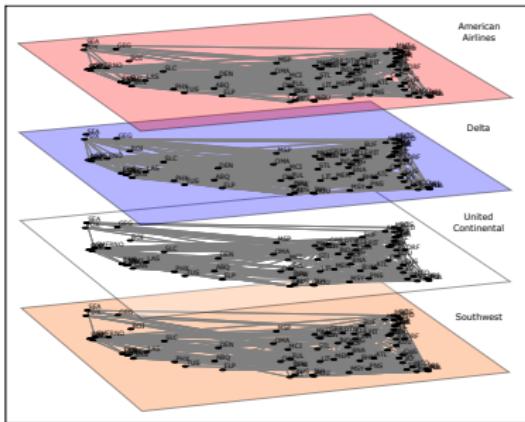
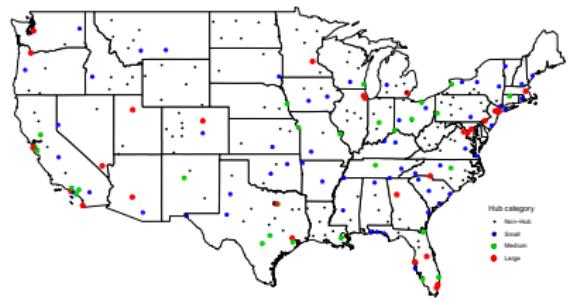
Demo video - finding injured after terrorist attack

https://www.youtube.com/watch?v=wyD005hF5tE&ab_channel=MattiasVillani

Airline network evolution

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

American data: 80 airports, 4 airlines, 40 quarters



Dynamic networks driven by latent variables

- Static Bernoulli model for adjacency matrix \mathbb{Y}

$$Y_{uv}(t) \stackrel{\text{iid}}{\sim} \text{Bern}(\pi)$$

- Dynamic Bernoulli with global latent GP

$$Y_{uv}(t) | \pi(t) \stackrel{\text{iid}}{\sim} \text{Bern}(\pi(t))$$

$$\text{Logit}[\pi(t)] = \log\left(\frac{\pi(t)}{1 - \pi(t)}\right) = z(t),$$

$$z(t) \sim \text{GP}(\mu(t), K(t, t'))$$

- Dynamic Bernoulli with latent GPs at nodes

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

$$\text{Logit}[\pi_{uv}(t)] = z(t) - d(x_u(t), x_v(t)),$$

$$z(t) \sim \text{GP}(\mu_z(t), K_z(t', t))$$

$$x_u(t) \sim \text{GP}(\mu_u(t), K_u(t', t)), u = 1, \dots, N.$$

Dynamic multi-layer networks

- Dynamic multi-layer Bernoulli model² for layer $k = 1, \dots, K$:

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

$$\text{Logit}\left[\pi_{uv}^{(k)}(t)\right] = z(t) - d(x_u(t), x_v(t)) - d\left(x_u^{(k)}(t), x_v^{(k)}(t)\right),$$

where

- ▶ Global GP across airport and airlines: $z(t)$
 - ▶ Airport-specific GPs, but shared across airlines: $x_u(t)$.
 - ▶ Airport/Airline-specific GPs: $x_u^{(k)}(t)$.
- Number of GPs: $1 + N + KN$ for N nodes and K layers.
 - Regularization is important.

²Durante et al (2017). Bayesian learning of dynamic multilayer networks, JMLR.

Stochastic block multi-layer networks

■ Stochastic block model for multi-layer networks³

$$Y_{uv}^{(k)}(t) | (s_u = a, s_v = b) \sim \text{Bern}\left(\pi_{ab}^{(k)}(t)\right)$$

$$\text{Logit}\left[\pi_{ab}^{(k)}(t)\right] = z(t) - d(x_a(t), x_b(t)) - d\left(x_a^{(k)}(t), x_b^{(k)}(t)\right),$$

$$\text{Logit}\left[\pi_{aa}^{(k)}(t)\right] = z(t) + \beta x_a(t) + \gamma x_a^{(k)}(t),$$

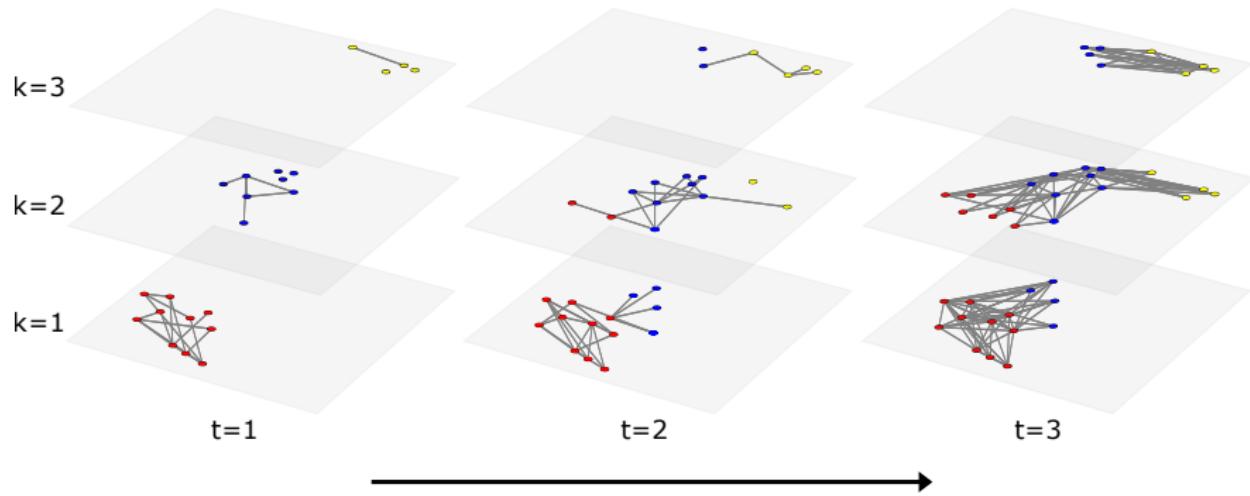
$$s_u \sim \text{Categorical}(\omega_1, \dots, \omega_B)$$

- Number of GPs: $1 + B + KB$ for B blocks. $B \ll N$.
- Need to learn the latent block indicators, s_u , for $u = 1, \dots, N$.
- Gibbs sampling via Pólya-Gamma augmentation.

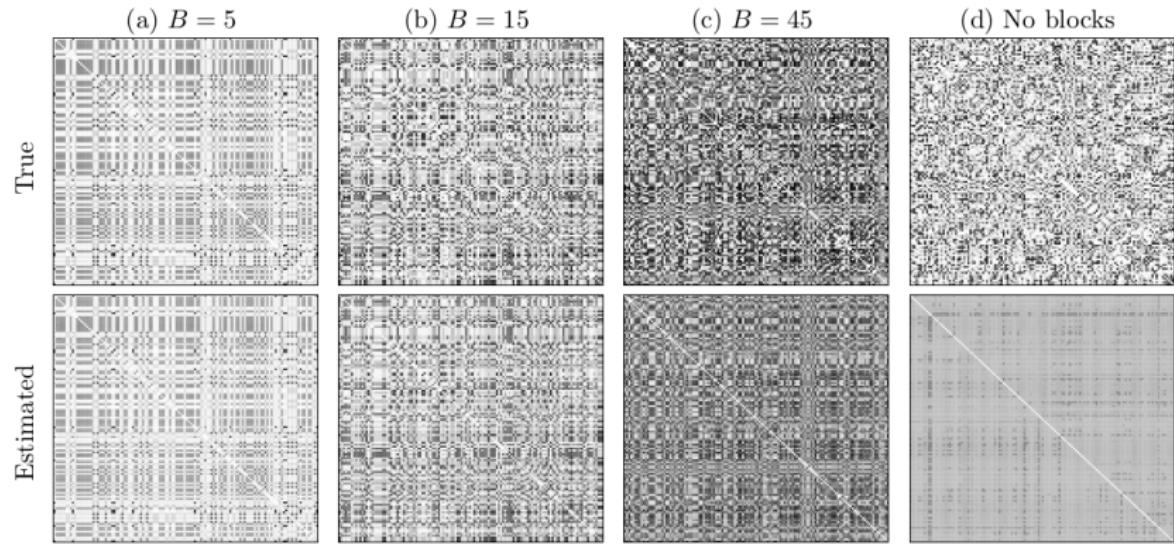
³Rodriguez-Deniz et (2022). A Multilayered Block Network Model to Forecast Large Dynamic

Transportation Graphs: an Application to US Air Transport, *Transportation Research C*.

Stochastic block multi-layer networks

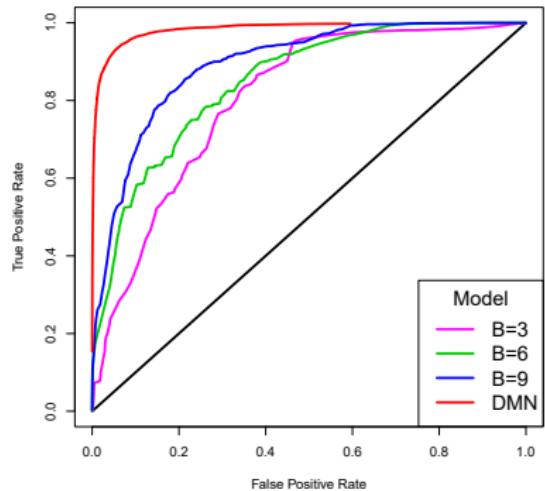


Fitting a $B = 10$ block multi-layer model

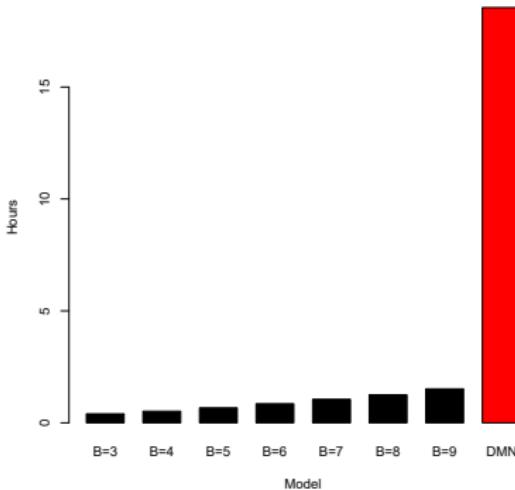


Predictive performance on 4 quarters test data

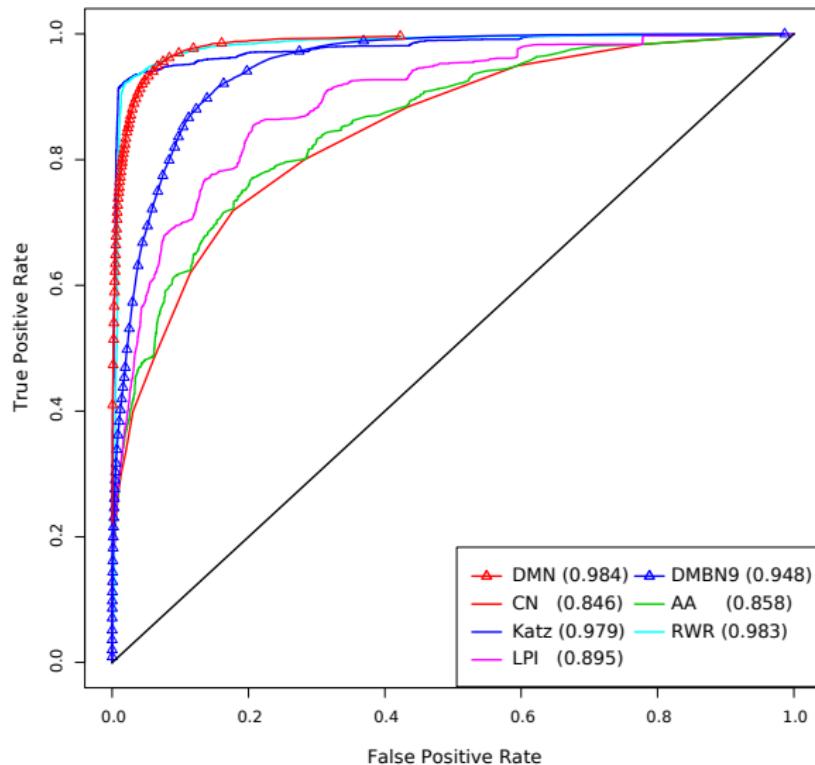
ROC curves



Computing times



Comparison to non-probabilistic SOTA



Bayesian optimization

- Aim: **maximization of expensive function**

$$\operatorname{argmax}_{x \in \mathcal{X}} f(x)$$

- **Bayesian optimization:**

- ▶ Assume $f \sim \mathcal{GP}$
- ▶ Evaluate f at x_1, x_2, \dots, x_n .
- ▶ Update to posterior distribution $f|f(x_1), \dots, f(x_n) \sim \mathcal{GP}$.
- ▶ Use posterior of f to find a new x_{n+1} . **Explore vs Exploit**.
- ▶ Iterate until convergence.

- Optimal x_{n+1} through an **acquisition function**.

Acquisition functions

- Probability of Improvement (PI)

$$a_{\text{PI}}(\mathbf{x}) \equiv \Pr(f(\mathbf{x}) > f_{\text{best}}),$$

where f_{best} is the highest function value so far.

- Expected Improvement (EI) takes also into account the size of the improvement.
- Non-convex acquisition function optimization, but deterministic and cheaper than original problem.
- Expected Improvement per Second - takes a known function cost into account.

Marginal likelihood estimated from sampling

- Marginal likelihood $f(\theta) \equiv \ln p(\mathbf{Y}_{1:T}|\theta)$ often estimated by Monte Carlo sampling.
- Noisy evaluations $\hat{f}(\theta)$.
- Precision of $\hat{f}(\theta)$ controlled via number of samples G .
- Sampling efficiency and therefore $\mathbb{V}(\hat{f}(\theta))$ varies over θ -space, particularly when θ contains prior/regularization hyperparameters.

Bayesian Optimization with Optimized Precision⁴

■ BOOP:

- ▶ Early stopping of evaluation when $\text{PI} < \alpha$.
- ▶ G random - we don't know G until we visit θ .
- ▶ EI per second, but with G predicted for every θ .

- Early stopping affects the planning of future computations.
- BOOP can try θ with lower EI, if expected to be cheap.
- Heteroscedastic GP regression model for the estimates

$$\hat{f}(\theta_i) = f(\theta_i) + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2(G_i))$$

■ GP for predicting the number of samples G :

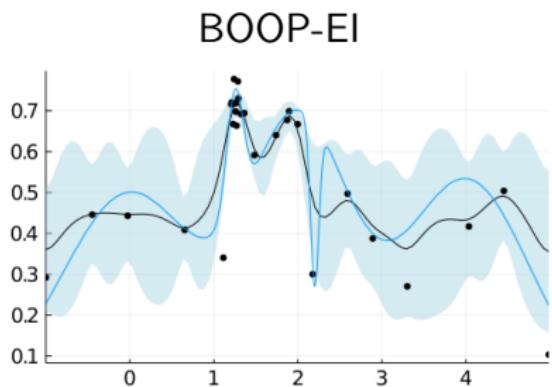
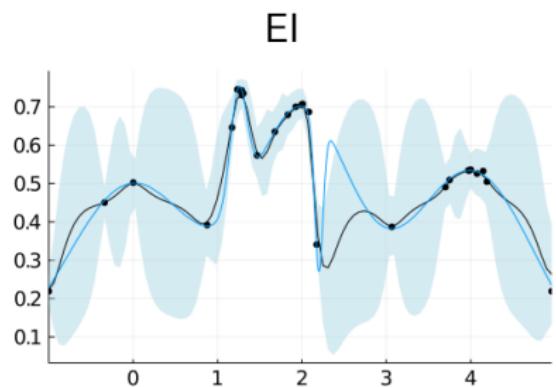
$$\ln G_i = h(z_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \psi^2),$$

where z are variables with predictive power for G , e.g. the hyperparameter values themselves, or $\hat{m}(\theta) - f_{\max}$.

⁴Gustafsson et (2022). Bayesian Optimization of Hyperparameters when the Marginal Likelihood is

Estimated by MCMC, arXiv.

BOOP in action - one example run



Bayesian VAR - BOOP finds estimates quicker

- 7 variable **steady-state BVAR** on US data.
- Gibbs sampling with **Chib's marginal likelihood estimator**.
- BO to find optimal prior hyperparameters:
 - Standard deviation first lag ($\lambda_1 \geq 0$)
 - Cross-equation shrinkage ($0 \leq \lambda_2 \leq 1$)
 - Lag length decay ($\lambda_3 \geq 0$)

