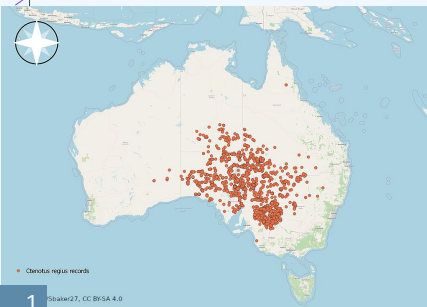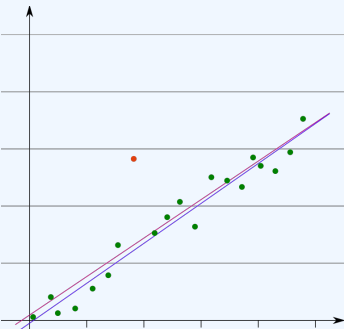# Gaussian processes for non-Gaussian likelihoods

ST John

Finnish Center for Artificial Intelligence
& Aalto University

Gaussian Process Summer School, 14 September 2021
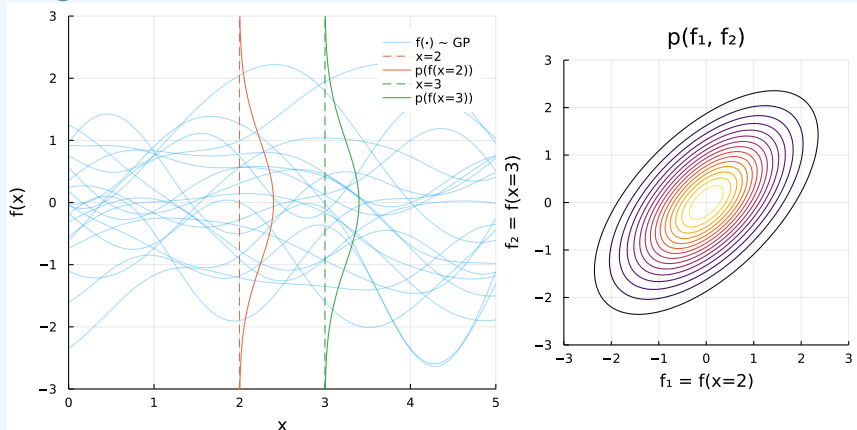
Chenotus regius records

Outline:

1. **Gaussian processes with Gaussian likelihood**
2. What is the likelihood? Connecting observations and Gaussian process prior
3. Non-Gaussian likelihoods: what happens to the posterior?
4. How to approximate the intractable
5. Comparisons

+ *Intuitive* understanding
+ Learning the language

– In-depth expertise
– Lots of maths

# SETTING THE SCENE

# GAUSSIAN PROCESS $f(\cdot)$

Distribution over *functions*
Marginals are Gaussian (mean and covariance)



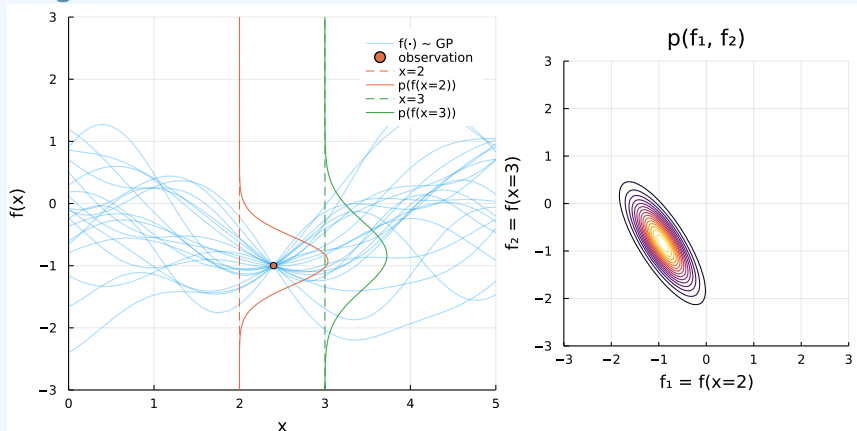`infinitecuriosity.org/vizgp`

Distribution over *functions*
Marginals are Gaussian (mean and covariance)



infinitecuriosity.org/vizgp

Without noise model, we interpolate observations:

$$y(x) = f(x)$$

Gaussian additive noise model, written two ways:

$$y(x) = f(x) + \epsilon, \qquad \epsilon \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2_{\text{noise}})$$
$$p(y \mid f) = \mathcal{N}(y \mid f, \sigma^2_{\text{noise}})$$

Gaussian additive noise model, written two ways:

$$y(x) = f(x) + \epsilon, \qquad \epsilon \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\text{noise}}^2)$$
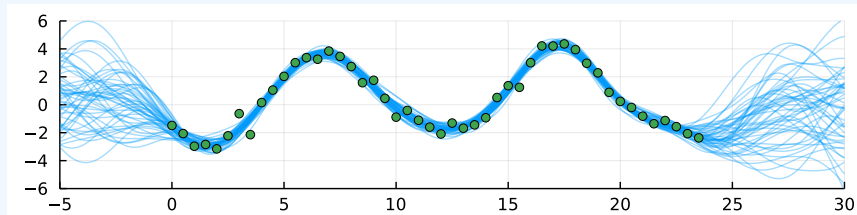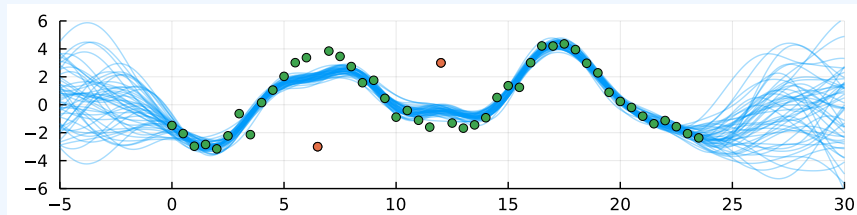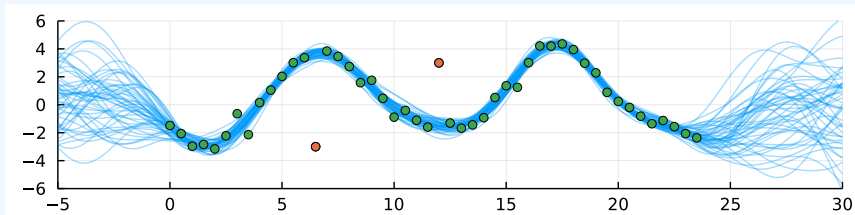$$p(y \mid f) = \mathcal{N}(y \mid f, \sigma_{\text{noise}}^2)$$
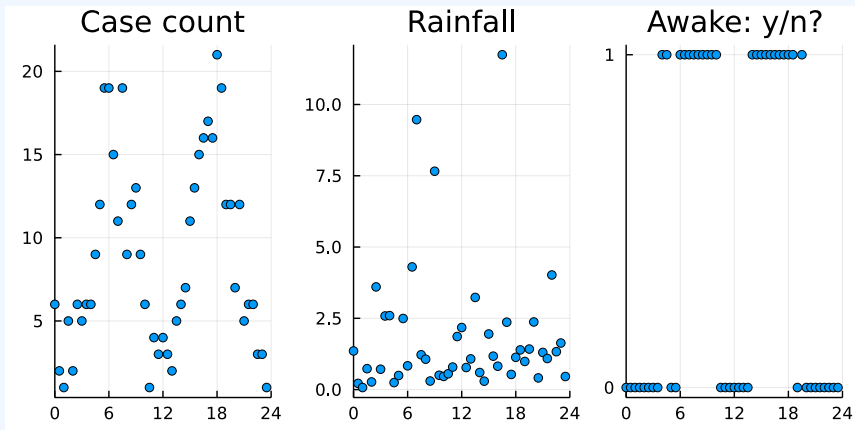
$$y(x) = f(x) + \epsilon, \qquad \epsilon \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\text{noise}}^2)$$
$$p(y \,|\, f) = \mathcal{N}(y \,|\, f, \sigma_{\text{noise}}^2)$$

# Likelihood

*latent* functional relationship

## Likelihood

$$p(\mathbf{y} \,|\, \mathbf{f}) = \prod_{i=1}^{N} p(y_i \,|\, f_i); \qquad f_i = f(x_i)$$
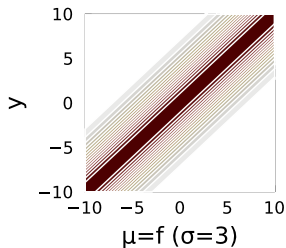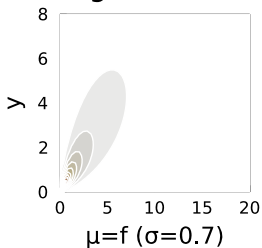
factorizing

$$p(y \,|\, f)$$

Function of two arguments:
$$y \mapsto p(y \,|\, f), \qquad f \mapsto p(y \,|\, f)$$
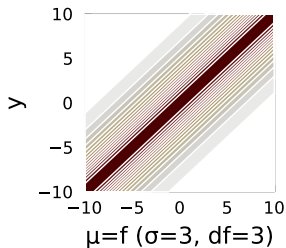
Two aspects of likelihoods:
   1. link functions
   2. log-concavity

$$\mathbb{E}[y] = \theta \in (0 \dots \infty) \qquad \text{link}(\theta) = f$$

$$f \sim \mathcal{N} \qquad \in (-\infty \dots \infty) \qquad \theta = \text{invlink}(f)$$

$\mathbb{E}[y] = \theta \in (0 \dots \infty)$

$f \sim \mathcal{N} \qquad \in (-\infty \dots \infty)$

$\mathrm{link}(\theta) = f$

$\theta = \mathrm{invlink}(f)$

Beta

Bernoulli

Poisson

$$\mathbb{E}[y] = \theta \in (0 \ldots \infty)$$

$$f \sim \mathcal{N} \qquad \in (-\infty \ldots \infty)$$

$$\text{link}(\theta) = f$$

$$\theta = \text{invlink}(f)$$

$$f(\alpha x + (1-\alpha)y) \geq \alpha f(x) + (1-\alpha)f(y)$$

# Posterior

## Likelihood

$$p(y \mid f)$$

## Joint distribution

$$p(y, f) = p(y \mid f)p(f)$$

## Posterior

$$f \mapsto p(f \mid y) = \frac{p(y \mid f)p(f)}{p(y)}$$

$$y \mapsto (f \mapsto p(f \mid y))$$

# Posterior predictions

At new point $x^*$:

$$p(f^* \mid x^*, \mathbf{x}, \mathbf{y}) = \int p(f^* \mid x^*, \mathbf{x}, \mathbf{f}) \, p(\mathbf{f} \mid \mathbf{x}, \mathbf{y}) \, \mathrm{d}\mathbf{f}$$

At training data:

$$p(\mathbf{f} \mid \mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{f} \mid \mathbf{x}) \prod_{i=1}^{N} p(y_i \mid f(x_i))}{\int p(\mathbf{f}' \mid \mathbf{x}) \prod_{i=1}^{N} p(y_i \mid f'(x_i)) \, \mathrm{d}\mathbf{f}''}$$
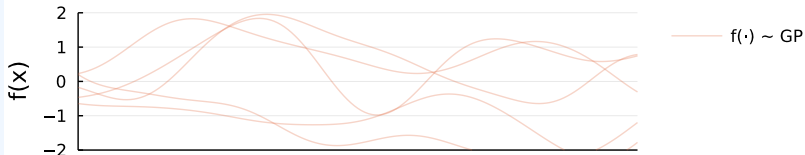
$$p(\mathbf{f} \mid \mathbf{y}) = \frac{1}{Z} p(\mathbf{f}) \prod_{i=1}^{N} p(y_i \mid f_i)$$

$$Z = p(\mathbf{y} \mid \mathcal{M}) = \int p(\mathbf{f} \mid \mathcal{M}) \prod_{i=1}^{N} p(y_i \mid f_i, \mathcal{M}) \, \mathrm{d}\mathbf{f}$$

"marginal likelihood" or "evidence" given **model** $\mathcal{M}$

$$p(\mathbf{f} \,|\, \mathbf{y}) = \frac{1}{Z} p(\mathbf{f}) \prod_{i=1}^{N} p(y_i \,|\, f_i)$$

Gaussian (process) prior $p(f(\cdot))$ …

  & Gaussian likelihood: conjugate case $\rightarrow$ posterior Gaussian

  & non-Gaussian $p(y|f) \rightarrow p(\mathbf{f} \,|\, \mathbf{y})$ also non-Gaussian, intractable

$$p(\mathbf{f}\,|\,\mathbf{y}) = \frac{p(\mathbf{f})\prod_{i=1}^{N}p(y_i\,|\,f_i)}{\int p(\mathbf{f'})\prod_{i=1}^{N}p(y_i\,|\,f_i')\mathrm{d}\mathbf{f'}}$$

$$f_1 = f(x_1)$$
$$f_2 = f(x_2)$$
$$\vdots$$
$$f_N = f(x_N)$$

- What is the likelihood $p(y \mid f)$?
- When is it non-Gaussian?
- Why does the posterior $p(f \mid y)$ become intractable?

Questions?! :)

# Approximations

- delta distribution
  - ▶ point estimate
- **Gaussian distribution**
  - ▶ Laplace
  - ▶ Expectation Propagation (EP)
  - ▶ Variational Bayes/Variational Inference (VB / VI)
- mixture of delta distributions
  - ▶ Markov Chain Monte Carlo (MCMC)
- mixture of Gaussians
- …



point estimate

Gaussian

delta mixture

Gaussian mixture

# GAUSSIAN APPROXIMATIONS

Approximating the posterior at observations:

$$p(\mathbf{f} \,|\, \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mu =?, \Sigma =?)$$

Predictions at new points:

$$p(f^* \,|\, x^*, \mathbf{y}) \approx q(f^*) = \int p(f^* \,|\, x^*, \mathbf{f})\, q(\mathbf{f})\, \mathrm{d}\mathbf{f}$$

# Demo: What does this mean for Gaussian processes?

tinyurl.com/nongaussian-inference-viz-v1

$$p(\mathbf{f} \mid \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mu =?, \Sigma =?)$$

match mean & variance at point

minimise divergence

**Laplace approximation**

Expectation Propagation (EP)

Variational Bayes (VB)

# LAPLACE APPROXIMATION

# Laplace approximation

Idea: log of Gaussian pdf = quadratic polynomial

$$p_{\mathcal{N}}(\mathbf{f}) = \frac{1}{\sqrt{(2\pi)^k|\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{f}-\mu)^{\top}\Sigma^{-1}(\mathbf{f}-\mu)\right)$$

Approximate quadratic polynomial:
2nd-order Taylor expansion of log of $h(f) = p(y\,|\,f)p(f)$ at $\hat{f}$

$$g(x+\delta) \approx g(x) + \left(\frac{\mathrm{d}g}{\mathrm{d}x}(x)\right)\delta + \frac{1}{2!}\left(\frac{\mathrm{d}^2g}{\mathrm{d}x^2}(x)\right)\delta^2$$

1. Find mode of posterior
   2nd-order gradient optimisation (e.g. Newton's method)
2. Match curvature (Hessian) at mode

$$p(f \mid y) = \frac{1}{Z} p(y \mid f) p(f)$$

$$\log p(f \mid y) = -\log Z + \log p(y \mid f) + \log p(f)$$

# $\log p(f \mid y) = -\log Z + \log h(f)$

$p(f \mid y) \approx \mathcal{N}(f \mid \hat{f}, -(\mathrm{d}^2 \log h / \mathrm{d}f^2)^{-1})$

marginal of 2D

# Laplace approximation: important properties

- find mode: Newton's method
- match curvature (Hessian) at mode
- "point estimate++"
- + simple, fast
- − poor approximation if mode is not representative (e.g. Bernoulli)
- − may not converge for non-log-concave likelihoods [3]

$$p(\mathbf{f} \mid \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mu =?, \Sigma =?)$$

match mean & variance at point

**minimise divergence**

Laplace approximation

Expectation Propagation (EP)

Variational Bayes (VB)

# Minimising divergences

"Relative entropy", "information gain" *from q to p*

$$D_{KL}(p\|q) = KL[p(x)\|q(x)] = \mathbb{E}_{x \sim p}\big[\log \frac{p(x)}{q(x)}\big] = \int p(x)\big[\log \frac{p(x)}{q(x)}\big]\mathrm{d}x$$

- non-symmetric: $KL[p\|q] \neq KL[q\|p]$
- positive: $KL \geq 0$ (Gibbs' inequality)
- minimum: $KL[p\|q] = 0 \Leftrightarrow q = p$.

# DEMO: KL BETWEEN TWO GAUSSIANS

tinyurl.com/nongaussian-inference-viz-v1

$$p(\mathbf{f} \mid \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f} \mid \mu =?, \Sigma =?)$$

1. $\min \text{KL}[p(\mathbf{f} \mid \mathbf{y}) \| q(\mathbf{f})]$: **Expectation Propagation**
2. $\min \text{KL}[q(\mathbf{f}) \| p(\mathbf{f} \mid \mathbf{y})]$: Variational Bayes

# Expectation Propagation (EP)

Exact posterior:

$$p(\mathbf{f}\,|\,\mathbf{y}) \propto p(\mathbf{f}) \prod_{i=1}^{N} p(y_i\,|\,f_i)$$

Approximate posterior:

$$q(\mathbf{f}) \propto p(\mathbf{f}) \prod_{i=1}^{N} t_i(f_i)$$

$$t_i = Z_i \mathcal{N}(f_i\,|\,\tilde{\mu}_i, \tilde{\sigma}_i^2)$$

Adding and subtracting natural (canonical) parameters

$$\text{``}\min \text{KL}[p(\mathbf{f}\,|\,\mathbf{y})\|q(\mathbf{f})]\text{''} \qquad q(\mathbf{f}) \propto p(\mathbf{f})\prod_{i=1}^{N}\underbrace{t_i(f_i)}_{\text{site}\,\propto\,\mathcal{N}(f_i)}$$

For each site $i$:

1. marginalize $\int q(\mathbf{f})\,\mathrm{d}f_{j\neq i} = q(f_i) \quad \not\propto t_i(f_i)$

2. improve local approximation: $\min \text{KL}[q(f_i)\frac{p(y_i\,|\,f_i)}{t_i(f_i)}\|q(f_i)\frac{t_i'(f_i)}{t_i(f_i)}]$

   2.1 *cavity* distribution $q_{-i}(f_i) = \frac{q(f_i)}{t_i(f_i)} \quad \Leftrightarrow \quad q(f_i) = q_{-i}(f_i)t_i(f_i)$

   2.2 *tilted* distribution $q_{\setminus i}(f_i) = q_{-i}(f_i)p(y_i\,|\,f_i)$

   2.3 $\operatorname{argmin} \text{KL}[q_{-i}(f_i)p(y_i\,|\,f_i)\|\hat{q}]$ by moment-matching

   2.4 update site: $t_i'(f_i) = \frac{\hat{q}}{q_{-i}(f_i)} \quad \Leftrightarrow \quad \hat{q} = q_{-i}(f_i)\,t_i'(f_i)$

3. compute new $q'(\mathbf{f})$ (rank-1 update)

iteration 1

iteration 2

# DEMO: EP IN 2D

tinyurl.com/nongaussian-inference-viz-v1

marginal of 2D

- prior
- exact posterior
- Laplace approximation
- Expectation Propagation

- multiple passes required to converge
- moment-matching (e.g. covering multiple modes)
+ effective for classification
– not guaranteed to converge
– updates may be invalid (non-log-concave likelihoods)

$$p(\mathbf{f} \mid \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f} \mid \mu =?, \Sigma =?)$$

✓ min KL$[p(\mathbf{f} \mid \mathbf{y})\|q(\mathbf{f})]$: Expectation Propagation
2. min KL$[q(\mathbf{f})\|p(\mathbf{f} \mid \mathbf{y})]$: **Variational Bayes**

# Variational Bayes (VB)
# Variational Inference (VI)

# Variational Bayes (VB)

Idea:
minimise divergence between $p(f \mid y)$ and $q(f)$ the "other" way

$$\underset{\mu,\Sigma}{\arg\min} \; \text{KL}\left[q(f) \| p(f \mid y)\right]$$

$$KL[q(f)\|p(f|y)]$$

$$= \int q(f)\left[\log\frac{q(f)}{p(f\,|\,y)}\right]\mathrm{d}f = \int q(f)\left[\log q(f) - \log p(f\,|\,y)\right]\mathrm{d}f$$

$$= \int q(f)\left[\log q(f) - \log p(f) - \log p(y\,|\,f) + \log p(y)\right]\mathrm{d}f$$

$$= \int q(f)\left[\log\frac{q(f)}{p(f)}\right]\mathrm{d}f - \int q(f)\left[\log p(y\,|\,f)\right]\mathrm{d}f + \log p(y)$$

$$= KL[q(f)\|p(f)] - \int q(f)\left[\log p(y\,|\,f)\right]\mathrm{d}f + \log p(y)$$

$$\log p(y) = \int q(f)\left[\log p(y\,|\,f)\right]\mathrm{d}f - KL[q(f)\|p(f)] + KL[q(f)\|p(f|y)]$$

$$\log p(y) = \int q(f) \big[ \log p(y \,|\, f) \big] \mathrm{d}f - \mathrm{KL}[q(f)\|p(f)] + \mathrm{KL}[q(f)\|p(f|y)]$$

$$\geq \int q(f) \big[ \log p(y \,|\, f) \big] \mathrm{d}f - \mathrm{KL}[q(f)\|p(f)]$$

Lower bound on the (log-)evidence $p(y)$ = ELBO

Integral separates for a factorizing likelihood:

$$\int q(\mathbf{f})\big[\log p(\mathbf{y}\,|\,\mathbf{f})\big]\mathrm{d}\mathbf{f}$$
$$= \sum_{i=1}^{N} \int q(f_i)\big[\log p(y_i\,|\,f_i)\big]\mathrm{d}f_i$$

Evaluating the 1D integrals:

- analytic (e.g. Exponential, Gamma, Poisson)
- Gauss–Hermite quadrature
- Monte Carlo (e.g. multi-class classification)

marginal of 2D

prior
exact posterior
Laplace
EP
VB

- principled: directly minimising divergence from true posterior
- mode-seeking (e.g. multi-modal posterior: fits just one)
- + minimises a true lower bound $\rightarrow$ convergence
- – underestimates variance

# Markov Chain Monte Carlo

- Samples $x_1, \ldots, x_T$
- "Markov" = 1-step history
- $x_{t+1} \sim p(x_{t+1} \mid x_t)$, independent of $x_{t-1}, \ldots, x_1$

Generate samples $\{x_t\} \sim p(f \mid y)$

Requires:

- *unnormalized* posterior
  $h(f) = p(y \mid f)p(f)$
- Markov proposal $q(x' \mid x_t)$
- initial $x_0$



In each iteration $t$:

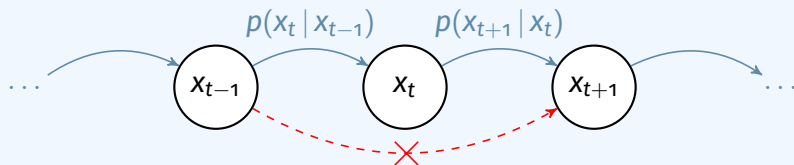1. Random proposal $x' \sim q(x' \mid x_t)$
2. Acceptance probability $\frac{h(x')}{h(x_t)} \rightarrow$ ensures sampling from $p(f \mid y)$

   accept: $x_{t+1} = x'$ \qquad reject: copy $x_{t+1} = x_t$

   $h(x') > h(x_t)$: always accepts $\rightarrow$ climbs uphill

# Demo: MCMC in 2D

tinyurl.com/nongaussian-inference-viz-v1

marginal of 2D

Legend: prior, exact posterior, Laplace, EP, VB, MCMC

- burn-in
- acceptance ratio
- auto-correlation, effective sample size (ESS); thinning to save memory
- mixing and multiple chains ($\hat{R}$)
- better proposals (HMC, NUTS) $\rightarrow$ use robust implementations
- + very accurate (gold-standard)
- – very slow, predictions require keeping all (thinned) samples around

Michael Betancourt's `betanalpha.github.io/writing/`

- Stan

- PyMC3

- TensorFlow Probability (GPflow)

- Turing.jl

✓ Gaussian processes with Gaussian likelihood
✓ What is the likelihood? Connecting observations and Gaussian process prior
✓ Non-Gaussian likelihoods: what happens to the posterior?
✓ How to approximate the intractable
    ✓ with Gaussians
        ■ Laplace
        ■ Expectation Propagation
        ■ Variational Bayes
    ✓ with samples: MCMC

5. **Comparisons**

# Comparison

**MCMC**

► samples

► gold standard

► slow

**Laplace**

► $\mathcal{N}$ = curvature at mode

► simple & fast

► often poor approximation

**EP**

► $\mathcal{N}$ matches marginal moments

► good calibration in classification

► may not converge

**Variational Bayes**

► $\mathcal{N}$ minimises $KL[q(f)\|p(f\,|\,y)]$

► principled, any likelihood

► underestimates variance
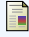
# What we did not cover...

- Marginal likelihood approximations for hyperparameter learning [6]
- How parametrisation affects Gaussianity of $p(f \mid y)$
- Connections between EP and VB ("PowerEP") [1]
- Combinations of MCMC and variational methods
- Augmenting likelihood with auxiliary variable $\rightarrow$ conditionally conjugate model [2]

# Questions!

📄 Thang D. Bui, Josiah Yan, and Richard E. Turner.
**A unifying framework for gaussian process pseudo-point approximations using power expectation propagation.**
*Journal of Machine Learning Research*, 18(104):1–72, 2017.

📄 Théo Galy-Fajou, Florian Wenzel, and Manfred Opper.
**Automated augmented conjugate inference for non-conjugate gaussian process models, 2020.**

📄 **Marcelo Hartmann and Jarno Vanhatalo.**
**Laplace approximation and natural gradient for gaussian process regression with heteroscedastic student-t model.**
*Statistics and Computing*, 29(4):753–773, October 2018.

📄 James Hensman, Nicolo Fusi, and Neil D. Lawrence.
**Gaussian processes for big data.**
*UAI*, 2013.

Pasi Jylänki, Jarno Vanhatalo, and Aki Vehtari.
**Robust gaussian process regression with a student-*t* likelihood.**
*Journal of Machine Learning Research*, 12(99):3227–3257, 2011.

Hannes Nickisch and Carl Edward Rasmussen.
**Approximations for binary gaussian process classification.**
*Journal of Machine Learning Research*, 9(67):2035–2078, 2008.