# Probabilistic Numerics
## Uncertainty in Computation

Philipp Hennig

Stuttgart
24 September 2018

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
**Intelligent Systems**

# The Numerics of Data Science & Machine Learning

nonlinear, non-analytic computations dominate the cost of data science

Machine Learning needs

Data —■— Model

**Computations**

| **integration** | | MCMC, VMP, EP, … | | probabilistic inference |
| **optimization** | like | SGD, Adam, RMSprop, … | for | (stochastic) fitting |
| **differential eqs.** | | Runge-Kutta, Multi-Step, … | | forecasting & control |
| **linear algebra** | | Cholesky, CG, spectral, … | | all of the above |

**generic methods** save design time, but do not address special needs

+ overly generic algorithms are inefficient
+ Big Data-specific challenges not addressed by "classical" methods

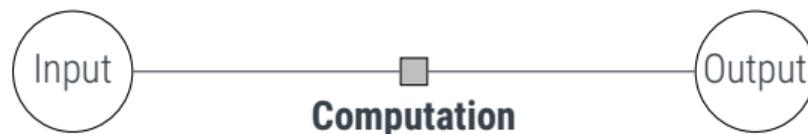Data Science / AI / ML needs to build its own numerical methods.

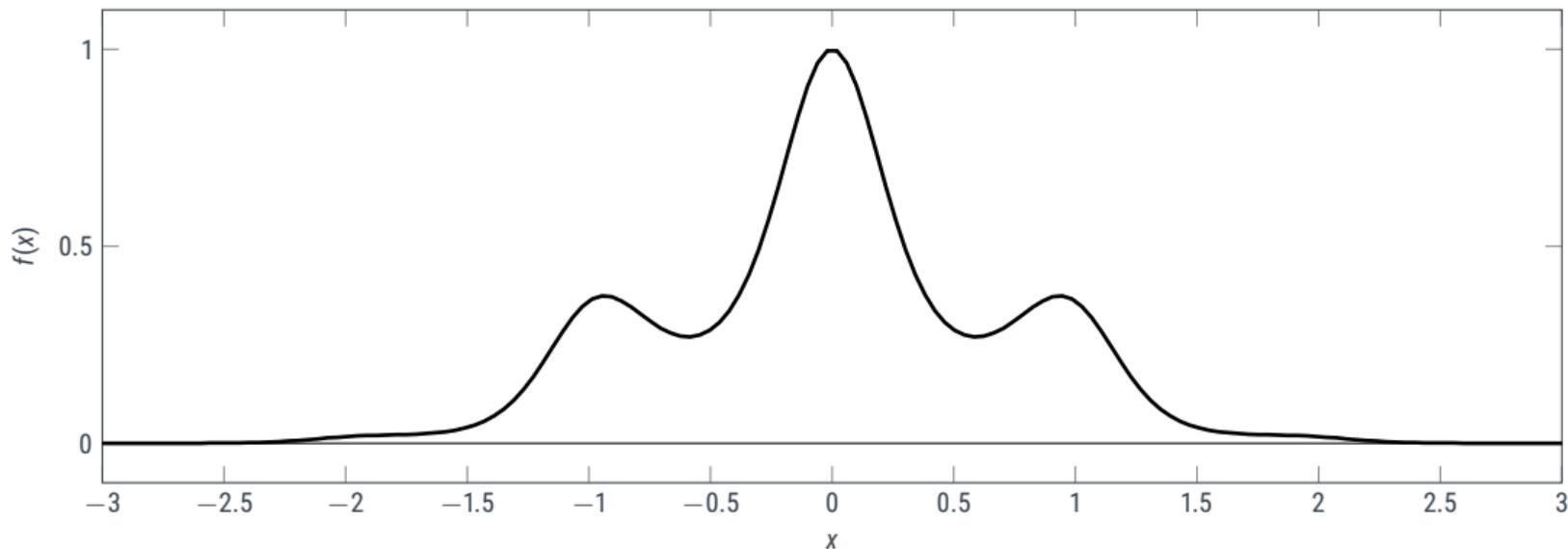As it turns out, we already have the right concepts!

Computation is Inference

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

http://probnum.org

[Poincaré 1896, Kimeldorf & Wahba 1970, Diaconis 1988, O'Hagan 1992, …]

Numerical methods **estimate latent** quantities **given** the result of computations.

| | | | |
|---|---|---|---|
| **integration** | estimate | $\int_a^b f(x)\,dx$ | given $\{f(x_i)\}$ |
| **linear algebra** | estimate | $x$ s.t. $Ax = b$ | given $\{As = y\}$ |
| **optimization** | estimate | $x$ s.t. $\nabla f(x) = 0$ | given $\{\nabla f(x_i)\}$ |
| **simulation** | estimate | $x(t)$ s.t. $x' = f(x, t)$ | given $\{f(x_i, t_i)\}$ |

It is thus possible to build
**probabilistic numerical methods**
that use **probability measures** as in- and outputs and assign **uncertainty** to computation.



Input ——————— ■ ——————— Output

**Computation**

2

# Integration
as Gaussian regression

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

$$f(x) = \exp(-\sin(3x)^2 - x^2) \qquad F = \int_{-3}^{3} f(x)\, dx = ?$$

A Wiener process prior $p(f, F)$…

Bayesian Quadrature

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[O'Hagan, 1985/1991]

$$p(f) = \mathcal{GP}(f; 0, k) \qquad k(x, x') = \min(x, x') + c$$

$$\Rightarrow p\left(\int_a^b f(x)\, dx\right) = \mathcal{N}\left[\int_a^b f(x)\, dx; \int_a^b m(x)\, dx, \iint_a^b k(x, x')\, dx\, dx'\right]$$

$$= \mathcal{N}(F; 0, -1/6(b^3 - a^3) + 1/2[b^3 - 2a^2 b + a^3] - (b - a)^2 c)$$

$$x_t = \arg \min \left[ \mathrm{var}_{p(F|x_1,\ldots,x_{t-1})}(F) \right]$$

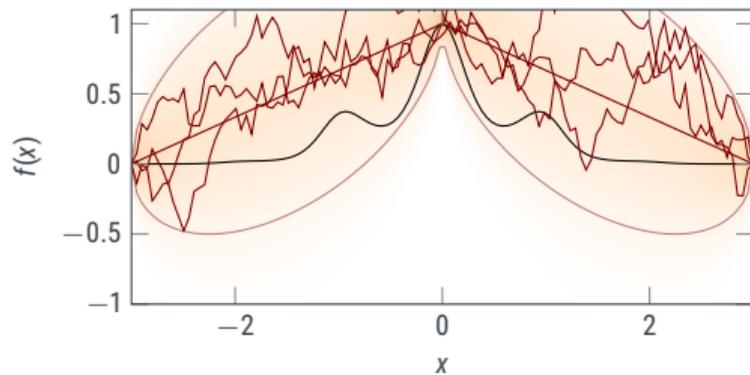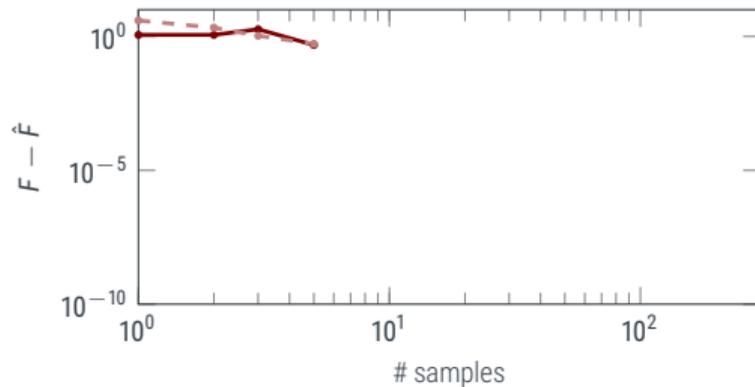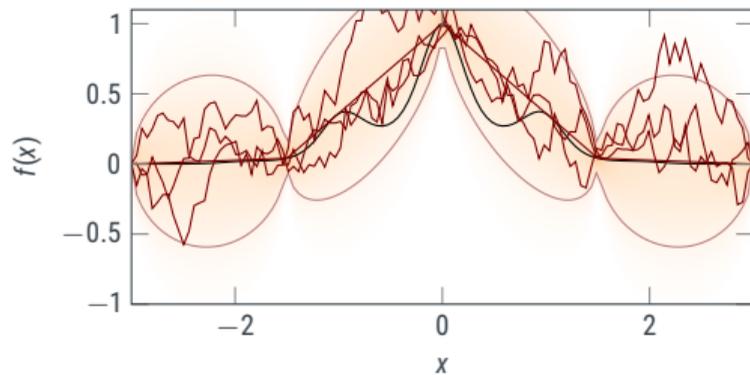+ **maximal reduction of variance** yields **regular grid**

...conditioned on actively collected information ...

UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

computation as the collection of information

$$x_t = \arg \min \left[ \mathrm{var}_{p(F|x_1,\ldots,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

. . . conditioned on actively collected information . . .

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems
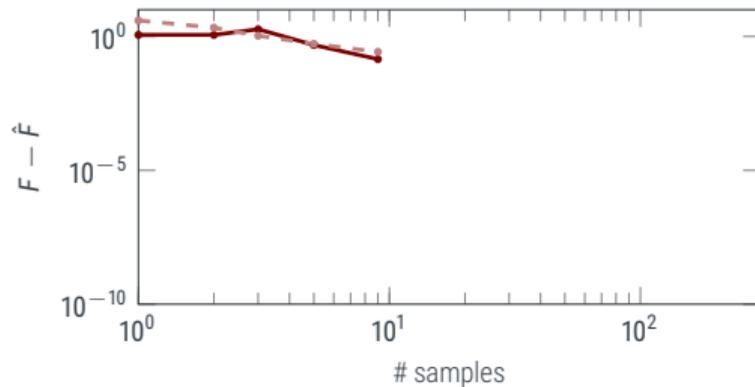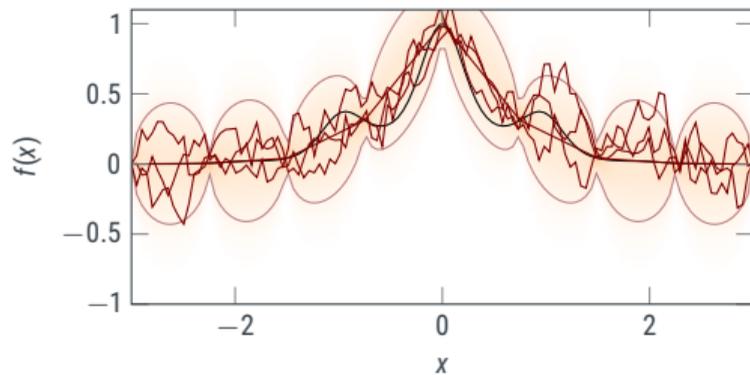
computation as the collection of information



$$x_t = \arg\min \left[ \operatorname{var}_{p(F|x_1,\dots,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

. . . conditioned on actively collected information . . .

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

computation as the collection of information

$$x_t = \arg \min \left[ \mathrm{var}_{p(F|x_1,...,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

$$x_t = \arg\min \left[ \mathrm{var}_{p(F|x_1,...,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

$$x_t = \arg\min \left[ \mathrm{var}_{p(F|x_1,\dots,x_{t-1})}(F) \right]$$

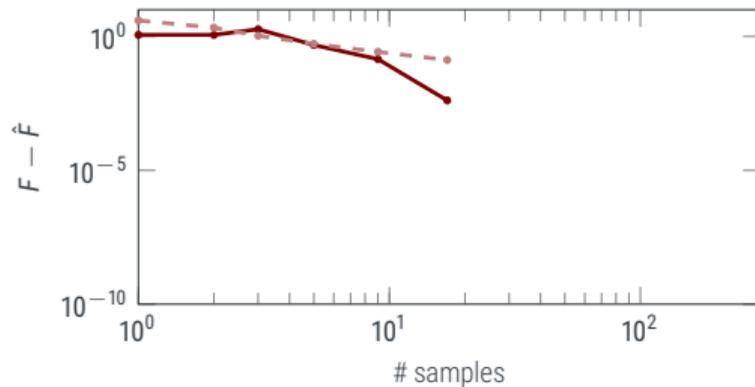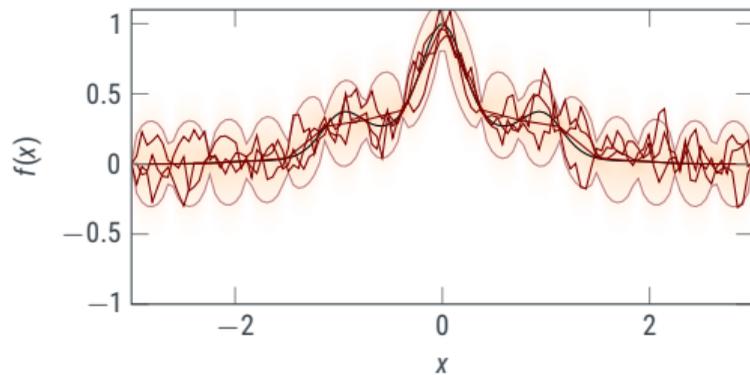+ **maximal reduction of variance** yields **regular grid**

$$x_t = \arg\min \left[ \text{var}_{p(F|x_1,\ldots,x_{t-1})}(F) \right]$$

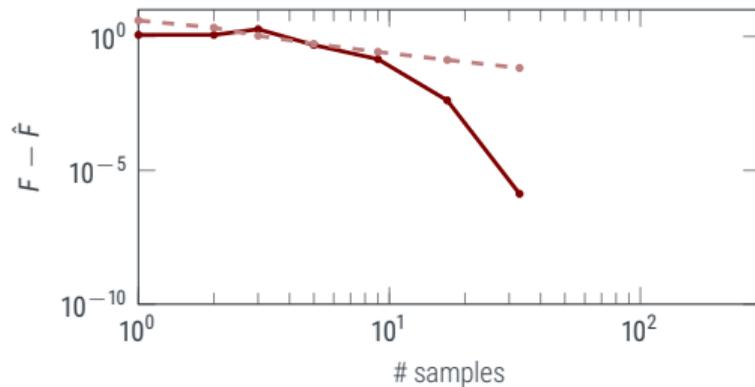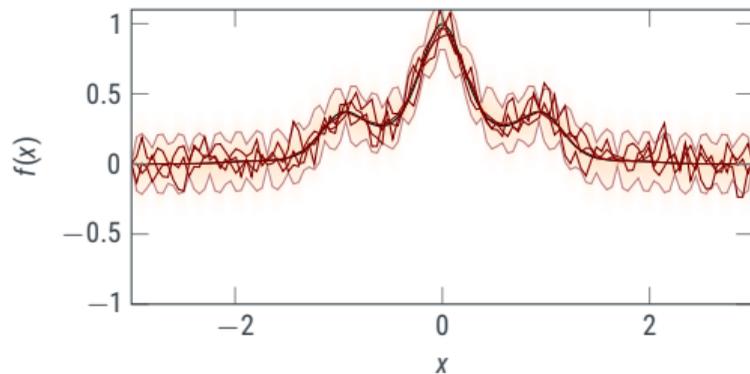+ **maximal reduction of variance** yields **regular grid**

... conditioned on actively collected information ...
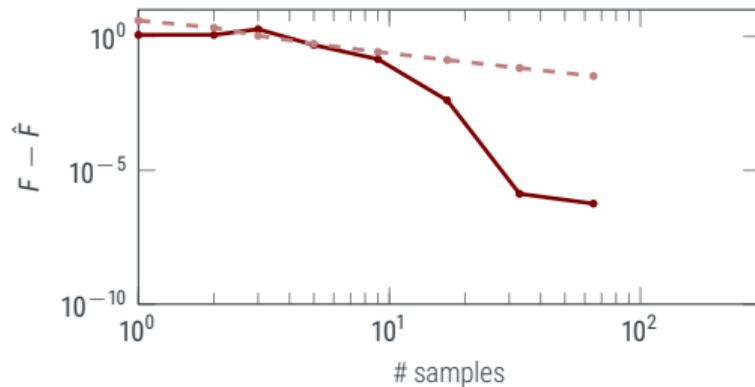
computation as the collection of information

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

$$x_t = \arg\min \left[ \text{var}_{p(F|x_1,...,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

…conditioned on actively collected information …

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

computation as the collection of information

$$x_t = \arg\min\left[\mathrm{var}_{p(F|x_1,\dots,x_{t-1})}(F)\right]$$

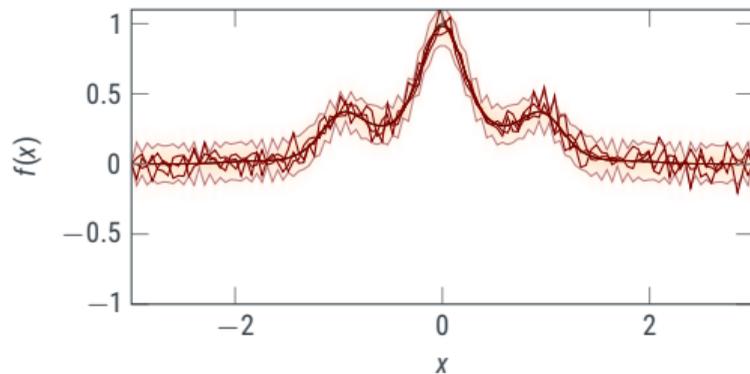+ **maximal reduction of variance** yields **regular grid**

$$x_t = \arg\min \left[ \text{var}_{p(F|x_1,\dots,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

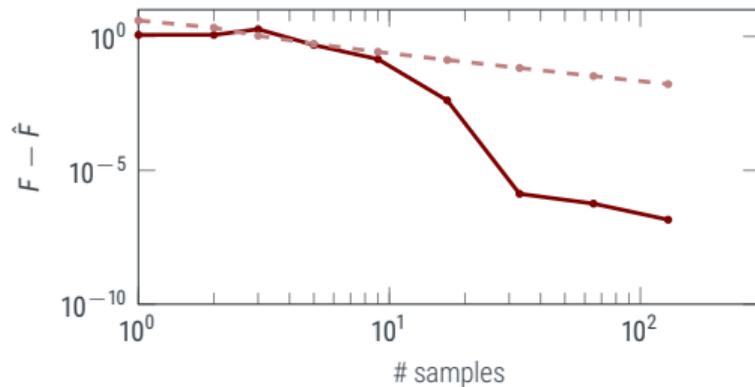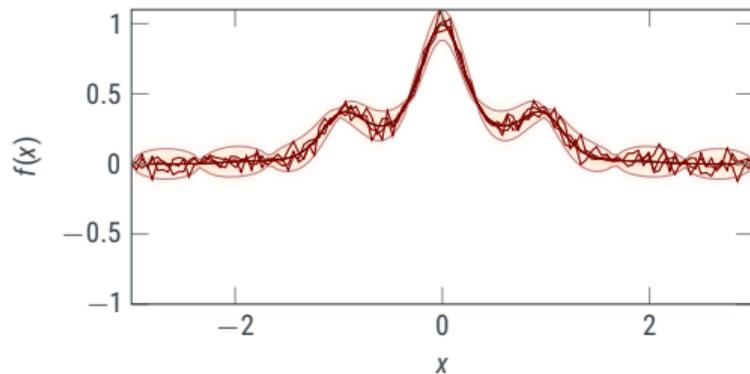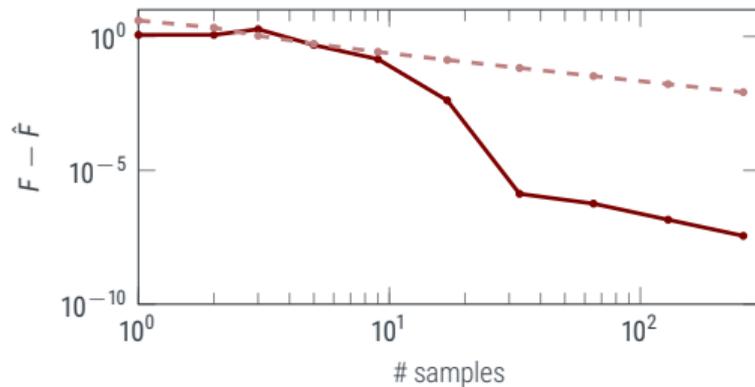$$x_t = \arg\min \left[ \mathrm{var}_{p(F|x_1,\ldots,x_{t-1})}(F) \right]$$

+ **maximal reduction of variance** yields **regular grid**

... yields the trapezoid rule!

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Kimeldorf & Wahba 1975, Diaconis 1988, O'Hagan 1985/1991]

$$E_{\mathbf{y}}[F] = \int E_{|\mathbf{y}}[f(x)] \, dx = \sum_{i=1}^{N-1} (x_{i+1} - x_i) \frac{1}{2} (f(x_{i+1}) + f(x_i))$$

+ **Trapezoid rule** is **MAP** estimate under Wiener process prior on *f*
+ regular grid is optimal expected information choice
+ error estimate is **under-confident**

# Computation as Inference

Bayesian inference on a latent (non-analytic) quantity from computable "observations"

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

Estimate *z* from computations *c*, under model *m*.

$$\underset{\text{posterior}}{p(z \mid c, m)} = \frac{\overset{\text{prior}}{p(z \mid m)}\,\overset{\text{likelihood}}{p(c \mid z, m)}}{\underset{\text{evidence}}{\int p(z \mid m)p(c \mid z, m)\,dz}}$$

# Classic methods as basic probabilistic inference

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

maximum a-posteriori estimation in Gaussian models

**Quadrature**

[Ajne & Dalenius 1960; Kimeldorf & Wahba 1975; Diaconis 1988; O'Hagan 1985/1991]

Gaussian Quadrature ⟷ GP Regression

---

**Linear Algebra**

[Hennig 2014]

Conjugate Gradients ⟷ Gaussian Regression

**Nonlinear Optimization**

[Hennig & Kiefel 2013]

BFGS / Quasi-Newton ⟷ Autoregressive Filtering

**Differential Equations**

[Schober, Duvenaud & Hennig 2014; Kersting & Hennig 2016; Schober & Hennig 2016]

Runge-Kutta; Nordsieck Methods ⟷ Gauss-Markov Filters

Probabilistic numerical methods can be as **fast** and **reliable** as classic ones.

Probabilistic ODE Solvers

Same story, different task

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$



There is a class of **solvers for initial value problems** that

+ has the same **complexity** as multi-step methods
+ has **high local approximation order** $q$ (like classic solvers)
+ has **calibrated posterior uncertainty** (order $q + 1/2$)
+ can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

Probabilistic ODE Solvers

Same story, different task

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$

There is a class of **solvers for initial value problems** that
  + has the same **complexity** as multi-step methods
  + has **high local approximation order $q$** (like classic solvers)
  + has **calibrated posterior uncertainty** (order $q + 1/2$)
  + can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

Probabilistic ODE Solvers

Same story, different task

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$

There is a class of **solvers for initial value problems** that
  + has the same **complexity** as multi-step methods
  + has **high local approximation order $q$** (like classic solvers)
  + has **calibrated posterior uncertainty** (order $q + 1/2$)
  + can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

# Probabilistic ODE Solvers

Same story, different task

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$

There is a class of **solvers for initial value problems** that

+ has the same **complexity** as multi-step methods
+ has **high local approximation order** $q$ (like classic solvers)
+ has **calibrated posterior uncertainty** (order $q + 1/2$)
+ can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

Probabilistic ODE Solvers

Same story, different task

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$

There is a class of **solvers for initial value problems** that
+ has the same **complexity** as multi-step methods
+ has **high local approximation order** $q$ (like classic solvers)
+ has **calibrated posterior uncertainty** (order $q + 1/2$)
+ can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

Probabilistic ODE Solvers

Same story, different task

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$

There is a class of **solvers for initial value problems** that
  + has the same **complexity** as multi-step methods
  + has **high local approximation order** $q$ (like classic solvers)
  + has **calibrated posterior uncertainty** (order $q + 1/2$)
  + can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

Probabilistic ODE Solvers

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

Same story, different task                                    [Schober, Duvenaud & P.H., 2014. Schober & P.H., 2016. Kersting & P.H., 2016]

$$x'(t) = f(x(t), t), \quad x(t_0) = x_0$$

There is a class of **solvers for initial value problems** that

+ has the same **complexity** as multi-step methods
+ has **high local approximation order** $q$ (like classic solvers)
+ has **calibrated posterior uncertainty** (order $q + 1/2$)
+ can use **uncertain initial value** $p(x_0) = \mathcal{N}(x_0; m_0, P_0)$

+ Computation is an instance of **inference**.
+ many classic numerical methods can be interpreted as probabilistic inference, arising from specific **generative models** (prior & likelihood)
+ Meaningful (calibrated) uncertainty can be constructed at minimal computational overhead (dominated by cost of point estimate)
+ **Designing a numerical method is a modelling task!**

The probabilistic viewpoint allows **new functionality** for **contemporary challenges**.

# New Functionality, and new Challenges

making use of the probabilistic numerics perspective

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

**Prior:** structural knowledge reduces complexity.

**Likelihood:**

$$p(z \mid c, m) = \frac{p(z \mid m) p(c \mid z, m)}{\int p(z \mid m) p(c \mid z, m) \, dz}$$

**Posterior:**

**Evidence:**

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

a prior specifically for integration of probability measures

+ $f > 0$ (*f* is probability measure)
+ $f \propto \exp(-x^2)$ (*f* is product of prior and likelihood terms)
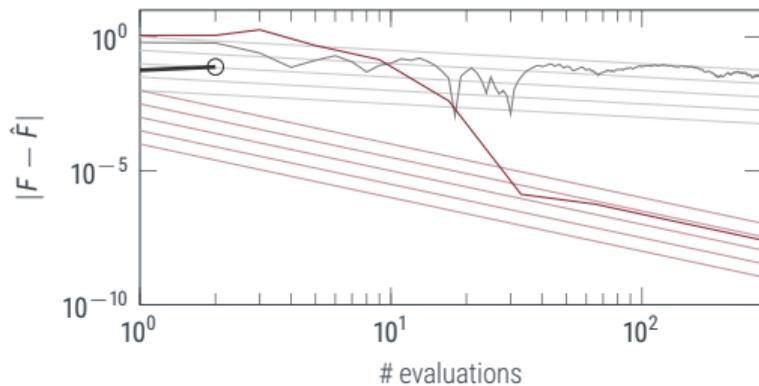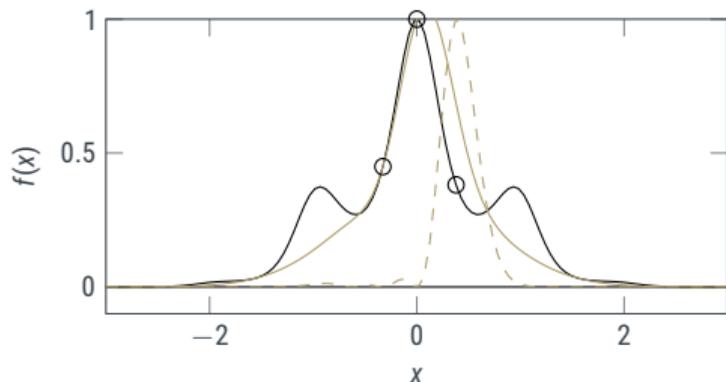+ $f \in \mathcal{C}^\infty$ (*f* is smooth)

Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
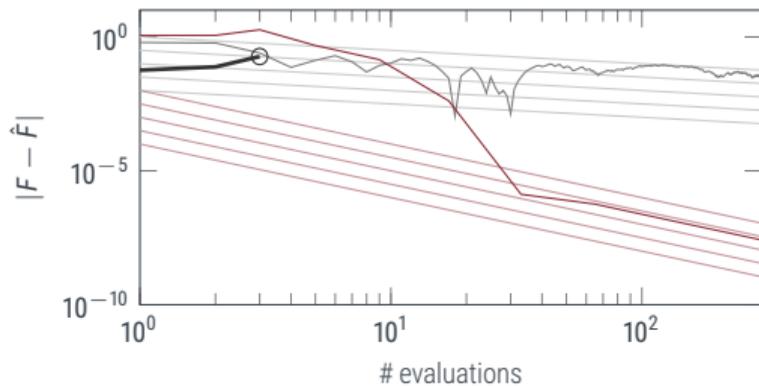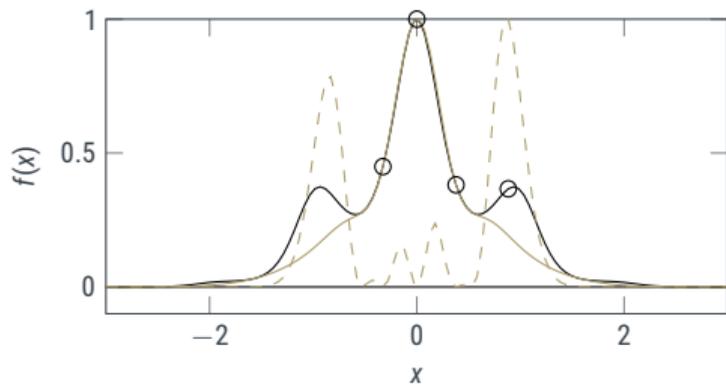+ faster (in wall-clock time) than MCMC

Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
+ faster (in wall-clock time) than MCMC

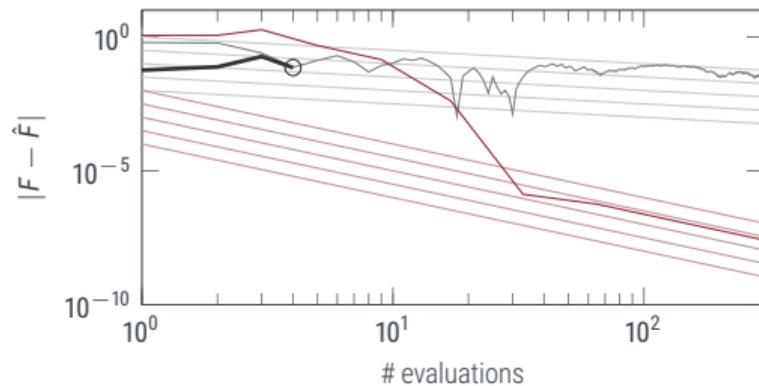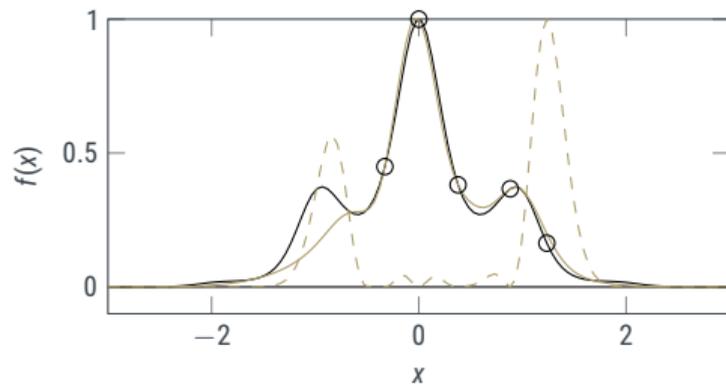Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
+ faster (in wall-clock time) than MCMC

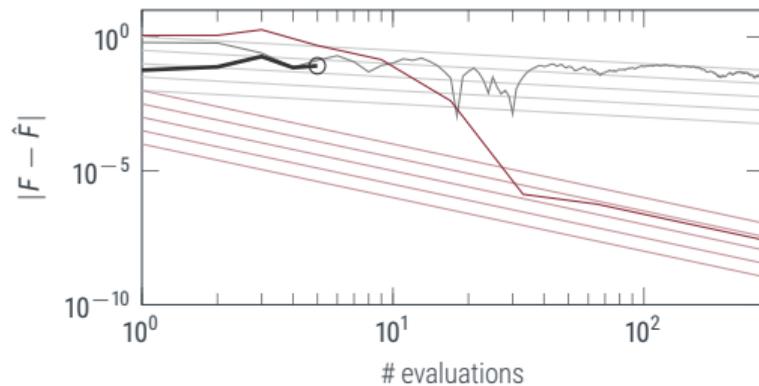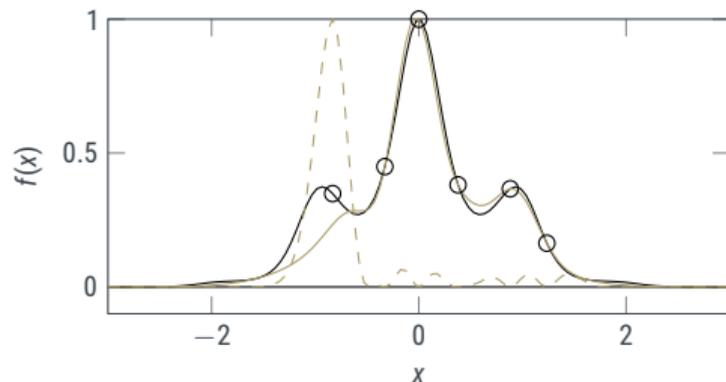Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
+ faster (in wall-clock time) than MCMC

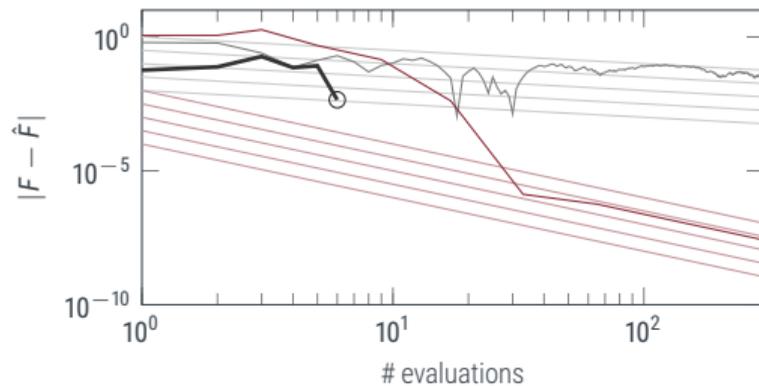Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
+ faster (in wall-clock time) than MCMC

Explicit prior knowledge yields reduces complexity.

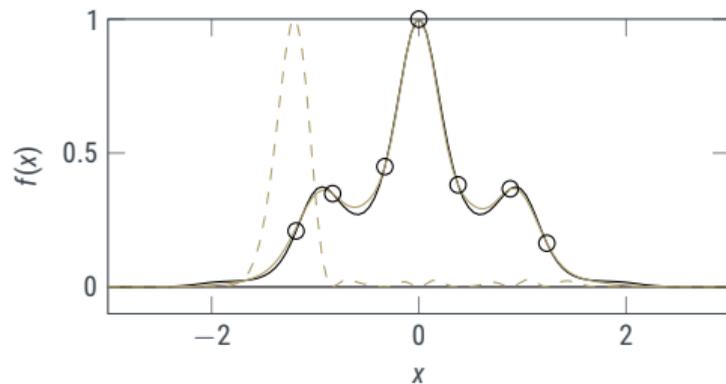[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
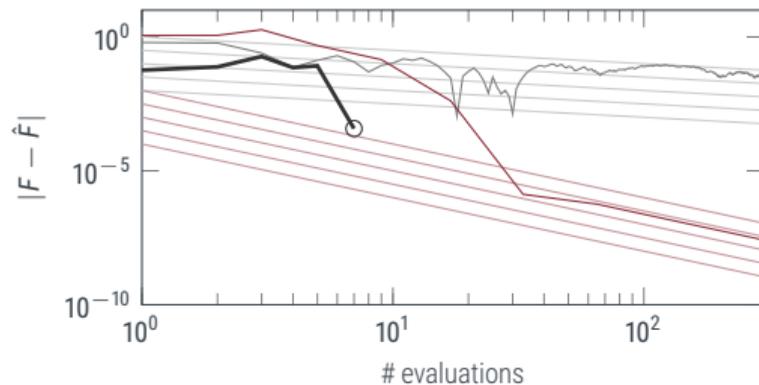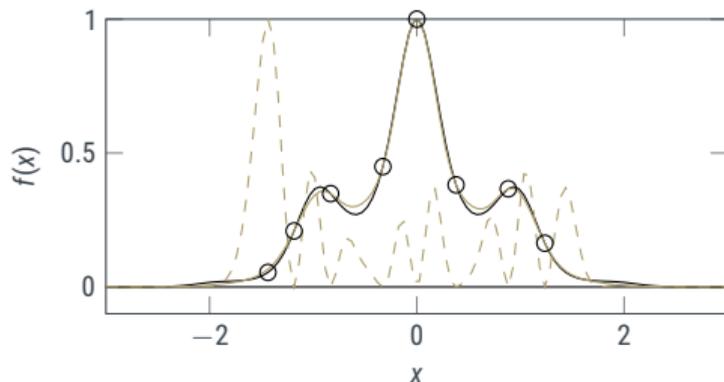+ faster (in wall-clock time) than MCMC

Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
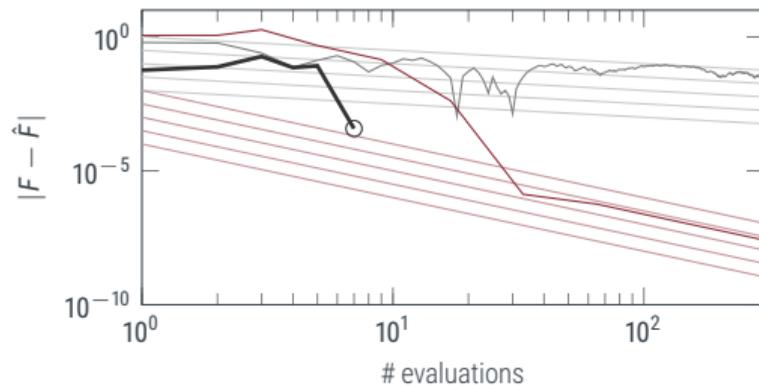+ faster (in wall-clock time) than MCMC

Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# An integration prior for probability measures

WArped Sequential Active Bayesian Integration (WSABI)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Gunter, Osborne, Garnett, Hennig, Roberts. NIPS 2014]

+ adaptive node placement
+ scales to, in principle, arbitrary dimensions
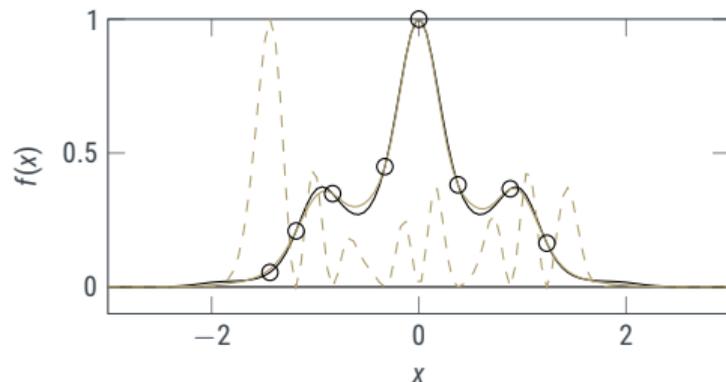+ faster (in wall-clock time) than MCMC

Explicit prior knowledge yields reduces complexity.

[cf. **information-based complexity**. E.g. Novak, 1988. Clancy et al. 2013, arXiv 1303.2412v2]

# Computation as Inference
new numerical functionality for machine learning

UNIVERSITÄT TÜBINGEN
EBERHARD KARLS

Intelligent Systems
Max Planck Institute for

Estimate *z* from computations *c*, under model *m*.

Prior: structural knowledge reduces complexity   Likelihood: modelling imprecision stabilizes algorithms

$$p(z \mid c, m) = \frac{p(z \mid m)p(c \mid z, m)}{\int p(z \mid m)p(c \mid z, m)\, dz}$$
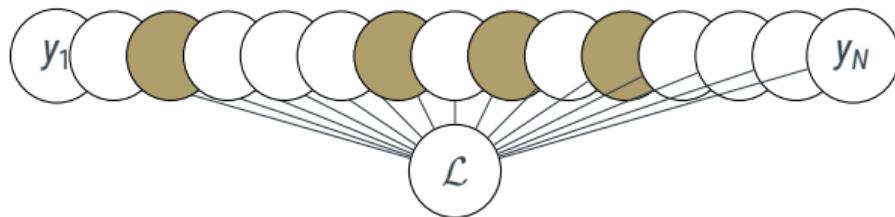
Posterior:   Evidence:

The usual assumption:

$$p(c \mid z, m) = \delta(c - A_m z)$$

In Big Data setting, iid. batching introduces Gaussian noise

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \ell(y_i; \theta) \approx \frac{1}{M} \sum_{j=1}^{M} \ell(y_j; \theta) =: \hat{\mathcal{L}}(\theta) \qquad M \ll N$$

$$p(\hat{\mathcal{L}} \mid \mathcal{L}) \approx \mathcal{N}\left(\hat{\mathcal{L}}; \mathcal{L}, \mathcal{O}\left(\frac{N-M}{NM}\right)\right)$$

# New numerics for Big Data
Uncertainty on Inputs directly effecting numerical decisions

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

In Big Data setting, iid. batching introduces Gaussian noise

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \ell(y_i; \theta) \approx \frac{1}{M} \sum_{j=1}^{M} \ell(y_j; \theta) =: \hat{\mathcal{L}}(\theta)$$

$$p(\hat{\mathcal{L}} \mid \mathcal{L}) \approx \mathcal{N}\left(\hat{\mathcal{L}}; \mathcal{L}, \mathcal{O}\left(\frac{N-M}{NM}\right)\right)$$

Contemporary machine learning requires tedious parameter fitting.

$$\theta_{t+1} = \theta_t - \alpha_t \nabla \hat{\mathcal{L}}(\theta_t)$$

+ step size / learning rate $\alpha_t$
+ batch size $M$
+ number of steps to termination
+ search directions



http://xkcd.com/1838

14

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \ell(y_i; \theta) \approx \frac{1}{M} \sum_{j=1}^{M} \ell(y_j; \theta) =: \hat{\mathcal{L}}(\theta) \qquad M \ll N \qquad p(\hat{\mathcal{L}} \mid \mathcal{L}) \approx \mathcal{N}\left(\hat{\mathcal{L}}; \mathcal{L}, \mathcal{O}\left(\frac{1}{M}\right)\right)$$
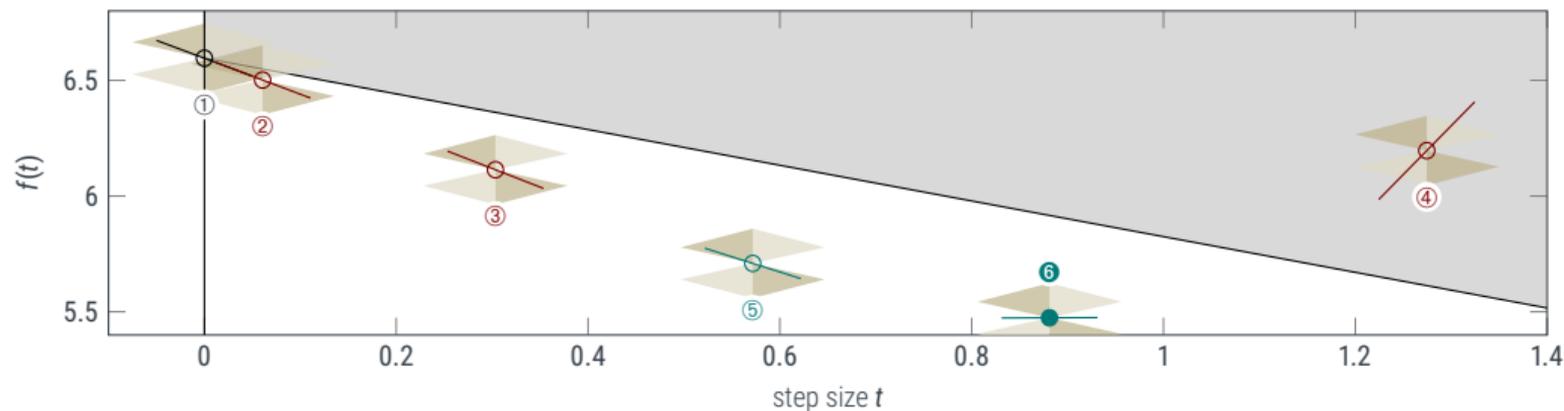
$$\text{var } \hat{\mathcal{L}}(\theta) \approx \frac{1}{M-1}\left(\frac{1}{M}\sum_{j=1}^{M}\ell^2(y_j;\theta) - \hat{L}^2(\theta)\right) \qquad p(\hat{\mathcal{L}} \mid \mathcal{L}) \approx \mathcal{N}\left(\hat{\mathcal{L}};\mathcal{L}, \text{var } \hat{\mathcal{L}}\right)$$

Capturing the likelihood requires a **new observable**! It's computation is not free, but cheap!
But without it, a key algorithmic parameter is **unidentified**!

# Choosing Step Sizes in the Presence of Noise

EBERHARD KARLS
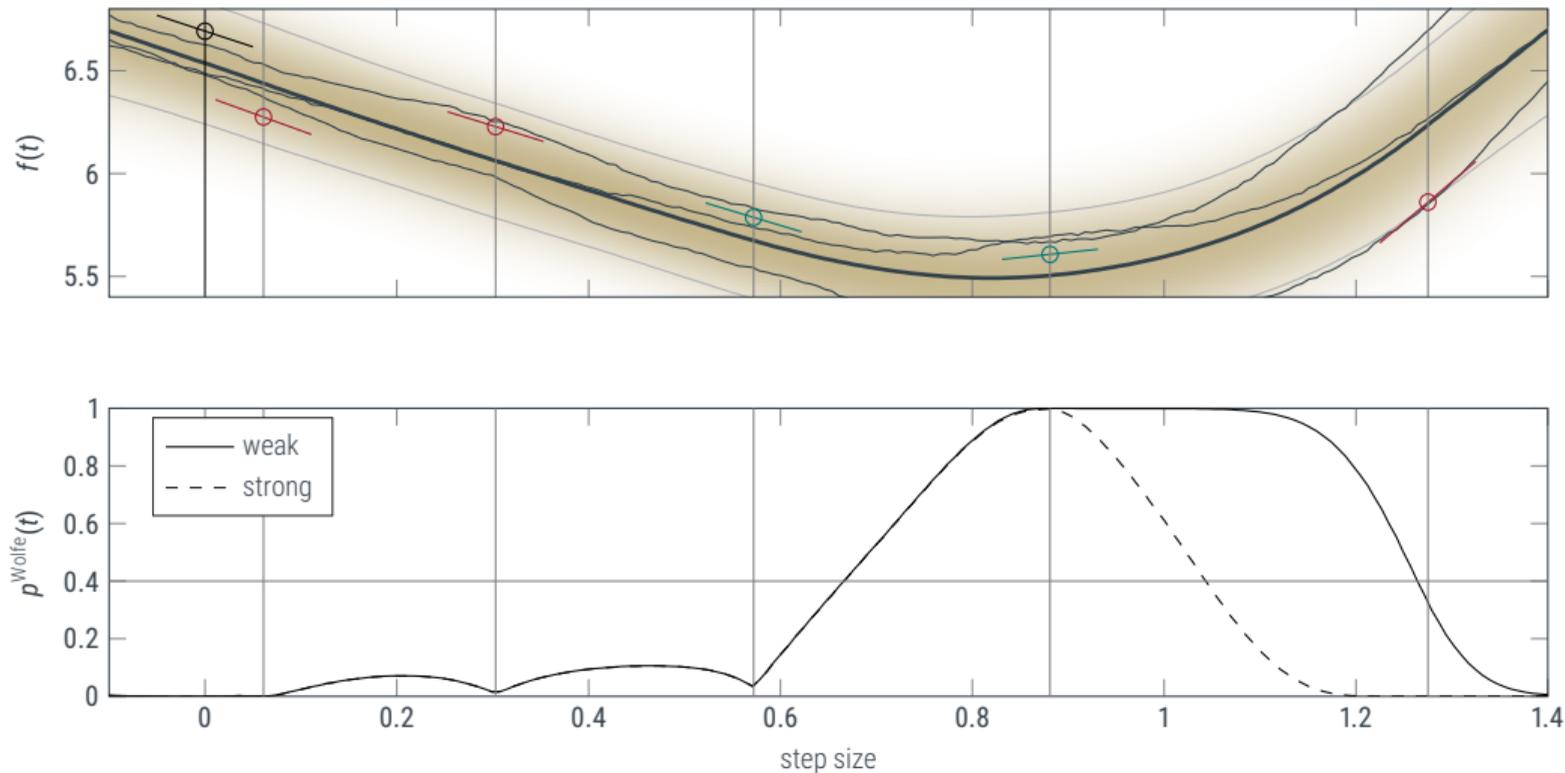UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

Probabilistic Line Searches

[Mahsereci & Hennig, NIPS 2015 (oral) / JMLR 2017]

+ $f'(t_{cand}) > 0$ ? bisect : extend
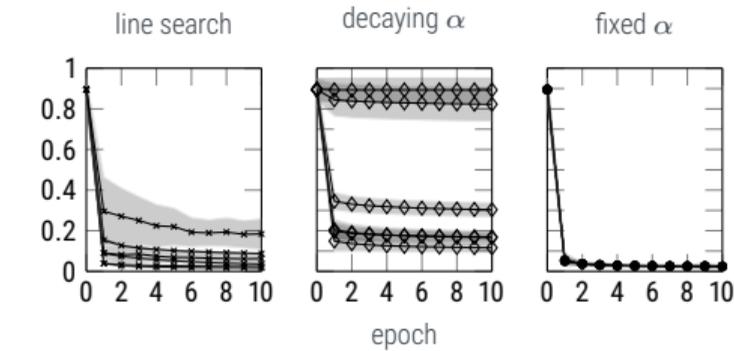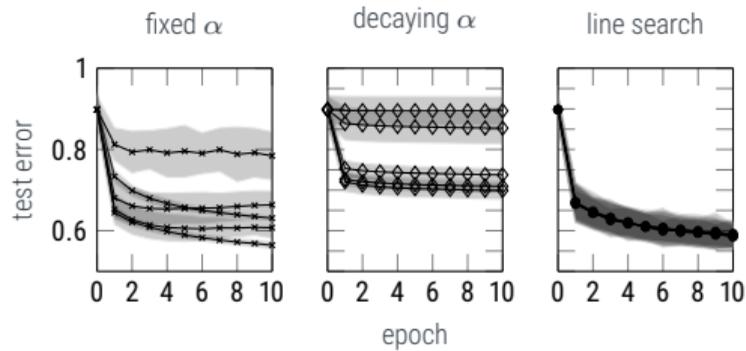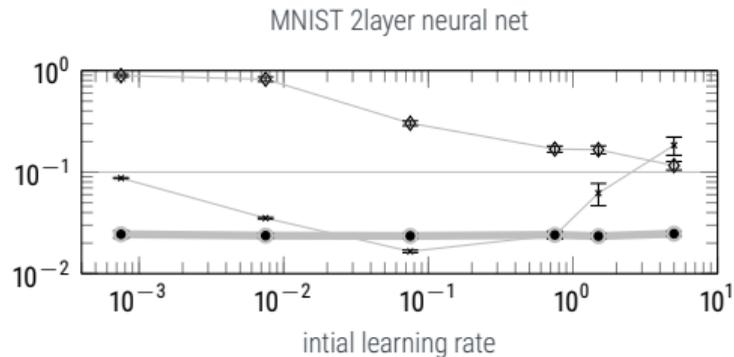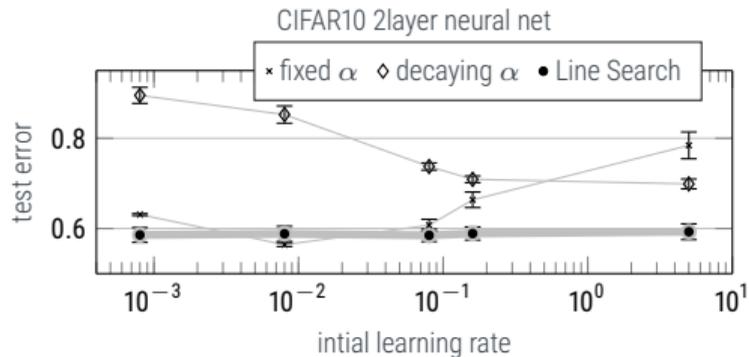+ until **Wolfe conditions** are fulfilled:

$$f(t) < f(0) + c_1 f'(0) \quad \text{AND} \quad |f'(t)| < c_2 |f'(0)|$$

# No more Learning Rates!

two-layer feed-forward perceptron. Details, additional results: Mahsereci & Hennig, JMLR **18**(119):1–59, 2017.

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

https://github.com/ProbabilisticNumerics/probabilistic_line_search

# Choosing Batch Sizes
trading off cost and precision

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Balles, Romero, Hennig, UAI 2017]

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \ell(y_i; \theta) \approx \frac{1}{M} \sum_{j=1}^{M} \ell(y_j; \theta) =: \hat{\mathcal{L}}(\theta) \qquad M \ll N$$

+ trade-off: $\text{std}[\nabla \hat{\mathcal{L}}] = \mathcal{O}(1/\sqrt{M})$, but cost is $\mathcal{O}(M)$

+ for SGD: lower bound on **improvement**: Assume $\nabla \mathcal{L}$ Lipschitz

$$\mathcal{L}(\theta_t) - \mathcal{L}(\theta_{t+1}) \geq G := \alpha \nabla \mathcal{L}(\theta_t)^\mathsf{T} \nabla \hat{\mathcal{L}}(\theta_t) - \frac{L\alpha^2}{2} \|\nabla \hat{\mathcal{L}}(\theta_t)\|^2$$

**expected improvement:** under $p(\hat{\mathcal{L}} \mid \mathcal{L})$ $\quad \mathbb{E}(G) = \left( \alpha - \frac{L\alpha^2}{2} \right) \|\nabla \mathcal{L}(\theta_t)\|^2 - \frac{L\alpha^2}{2M} \sum_\ell \text{var}[\nabla_\ell \hat{\mathcal{L}}(\theta_t)]$

+ maximize **expected improvement per cost**, let **line-search** find $\alpha = 1/L$, some further simplifications (local 2nd order approximation, assert $\min \mathcal{L} \gtrsim 0$),

$$M_* = \arg\max_M \frac{\mathbb{E}[G]}{M} \approx \alpha_t \frac{\sum_\ell \text{var}[\nabla_\ell \hat{\mathcal{L}}(\theta_t)]}{\hat{\mathcal{L}}(\theta_t)}$$

# Choosing Batch-Sizes

trading off cost and precision

EBERHARD KARLS UNIVERSITÄT TÜBINGEN

Max Planck Institute for Intelligent Systems



MNIST  SVHN  CIFAR-10  CIFAR-100

Train loss

Test accuracy

$M = 32$
$M = 128$
$M = 512$
adaptive

$M$

data read

https://github.com/ProbabilisticNumerics/CABS

20

Preventing Overfitting

early stopping without a validation set

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

[Mahsereci, Balles, Lassner, Hennig, arXiv 1703.09580]

+ in empirical risk minimization, just figuring out when to **stop** the optimizer is a non-trivial problem
+ even the full data set is a sample relative to the population
+ **overfitting** becomes a problem when gradients (with their estimatable variance) are statistically indistinguishable to white noise around zero

$$\log p(\nabla \hat{\mathcal{L}} \mid \nabla \mathcal{L} = 0) > \mathsf{E}_{p(\nabla \hat{\mathcal{L}} \mid \nabla \mathcal{L} = 0)} \left[ \log p(\nabla \hat{\mathcal{L}} \mid \nabla \mathcal{L} = 0) \right]$$

$$1 - \frac{M}{D} \sum_{\ell=1}^{D} \frac{(\nabla_\ell \mathcal{L}(\theta_t))^2}{\operatorname{var} \nabla_\ell \hat{\mathcal{L}}(\theta_t)} > 0 \quad \Rightarrow \quad \text{STOP!}$$

# Towards Black Box Deep Learning

inferring free parameters by hierarchical inference

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

+ **step sizes**
  *Probabilistic Line Searches for Stochastic Optimization*

  Mahsereci & Hennig
  NIPS 2015

  `https://github.com/ProbabilisticNumerics/probabilistic_line_search`

+ **batch sizes**
  *Coupling Adaptive Batch Sizes with Learning Rates*

  Balles, Romero, Hennig
  UAI 2017

  `https://github.com/ProbabilisticNumerics/cabs`

+ **termination criteria**
  *Early Stopping without a Validation Set*

  Mahsereci, Balles, Lassner, Hennig
  arXiv 1703.09580

+ **data sub-sampling** gives rise to imprecise computations / non-Dirac observations **likelihoods**
+ **free algorithmic parameters** may then become **un-identified**
+ likelihood shape can be identified with **minor computational overhead**
+ **classic methods** provide a **blue-print**
+ re-phrasing them probabilistically allow **inference** on free parameters

# Computation as Inference
new numerical functionality for machine learning

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

Estimate *z* from computations *c*, under model *m*.

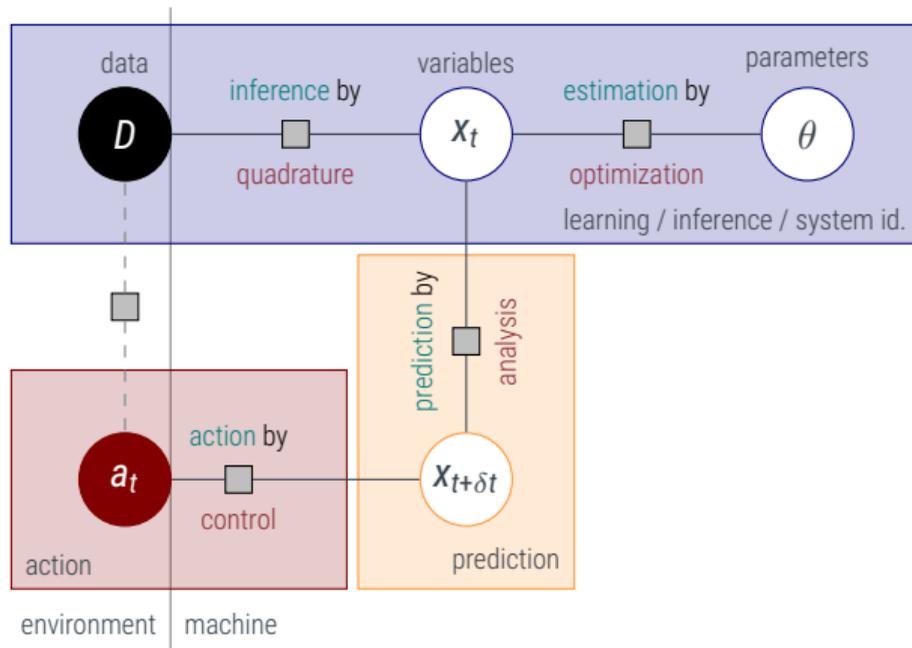**Prior:** structural knowledge reduces complexity    **Likelihood:** modelling imprecision stabilizes algorithms

$$p(z \mid c, m) = \frac{p(z \mid m)p(c \mid z, m)}{\int p(z \mid m)p(c \mid z, m)\, dz}$$

**Posterior:** tracking uncertainty for robustness    Evidence:

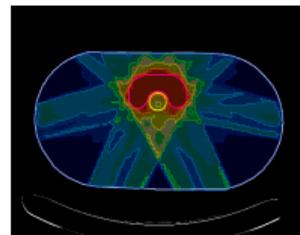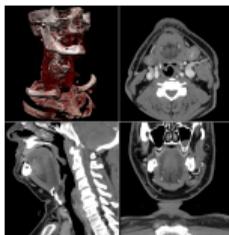cf. Hennig, Osborne, Girolami, Proc. Royal Soc. A, 2015

for some recent theory, see Thm. 5.9 in Cockayne, Oates, Sullivan, Girolami. arXiv 1702.03673
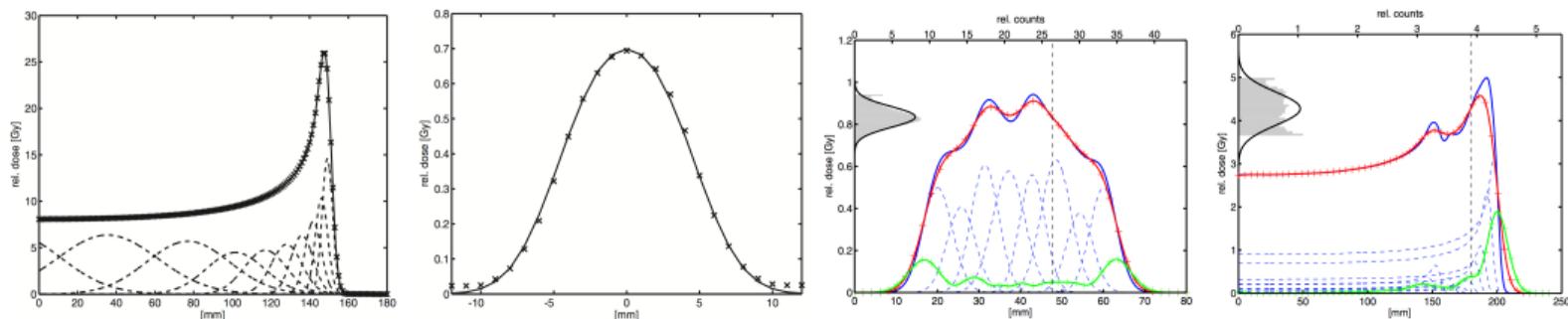
radiation treatment planning involves **approximately optimizing** an **imprecise** function subject to **uncertainties**.
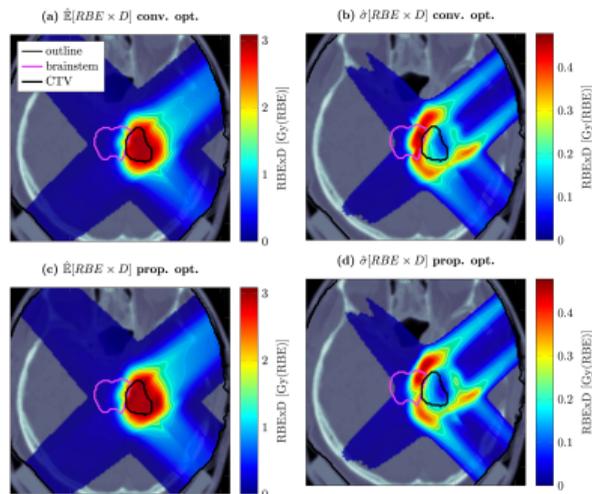
# Propagating Uncertainty through Pipelines

Analytical Probabilistic Treatment Planning — with DKFZ Heidelberg

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
**Intelligent Systems**

[Bangert et al., PMB, 2013, 2016, 2017]

+ map all involved non-linear functions into tractable (Hilbert-) space, with **quality guarantees**, bounds on approximation error

(a) $\hat{\mathbb{E}}[RBE \times D]$ conv. opt.    (b) $\hat{\sigma}[RBE \times D]$ conv. opt.

(c) $\hat{\mathbb{E}}[RBE \times D]$ prop. opt.    (d) $\hat{\sigma}[RBE \times D]$ prop. opt.

+ map all involved non-linear functions into tractable (Hilbert-) space, with **quality guarantees**, bounds on approximation error
+ track and **optimize uncertainties** across computation
+ to improve treatment outcome, reduce risk of complications

# Computation as Inference
new numerical functionality for machine learning

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

Estimate *z* from computations *c*, under model *m*.

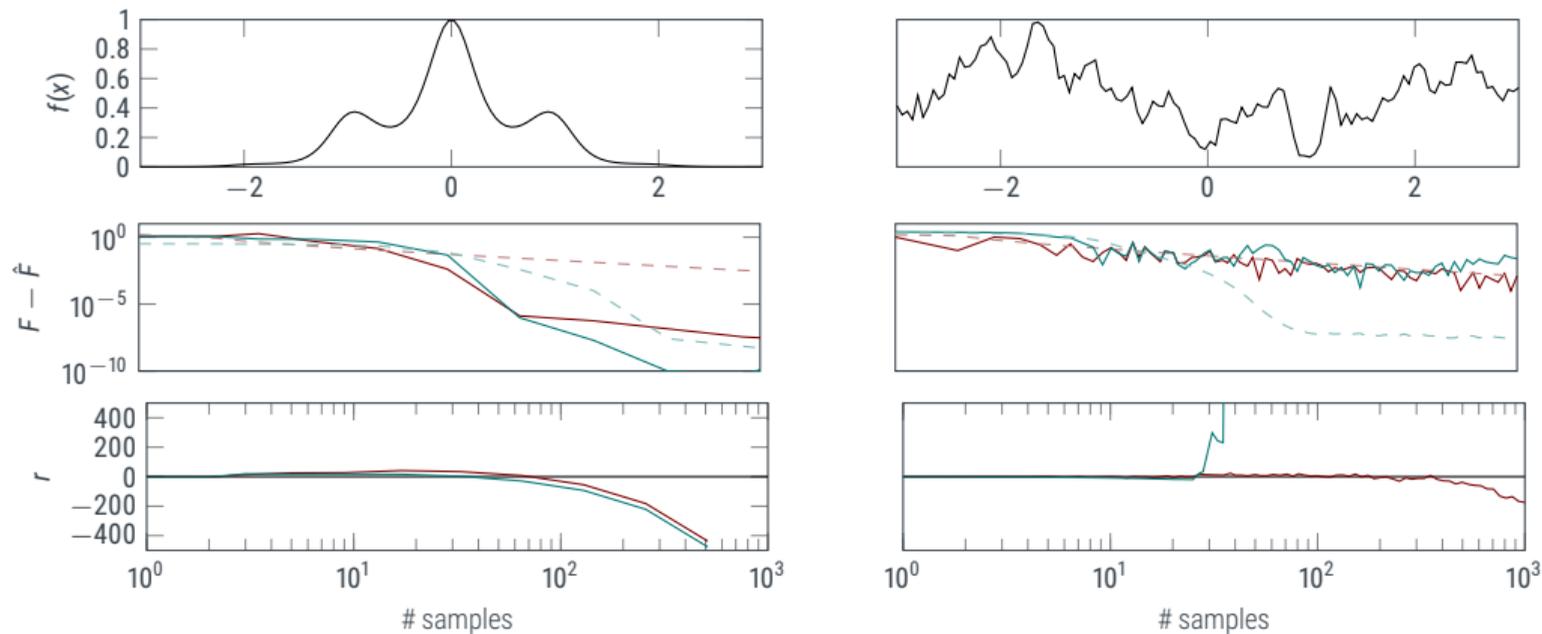**Prior:** structural knowledge reduces complexity

**Likelihood:** modelling imprecision stabilizes algorithms

$$p(z \mid c, m) = \frac{p(z \mid m)p(c \mid z, m)}{\int p(z \mid m)p(c \mid z, m)\, dz}$$

**Posterior:** tracking uncertainty for robustness

**Evidence:** checking models for safety

cf. Hennig, Osborne, Girolami, Proc. Royal Soc. A, 2015

$$r = \mathrm{E}_{\tilde{f}}\left[\log\frac{p(\tilde{f}(\mathbf{x}))}{p(f(\mathbf{x}))}\right] = (f(\mathbf{x}) - \mu(\mathbf{x}))^{\mathsf{T}} K^{-1}(f(\mathbf{x}) - \mu(\mathbf{x})) - N$$

# Summary

Uncertain computation **as** and **for** statistical modelling and machine learning

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Max Planck Institute for
Intelligent Systems

+ **computation is inference → probabilistic numerical methods**
  + probability measures for **uncertain** inputs and outputs
  + classic methods as special cases

Building numerical methods for contemporary challenges amounts to designing probabilistic models.

**prior:** structural knowledge reduces complexity

**likelihood:** imprecise computation lowers cost

**posterior:** uncertainty can be propagated through computations

**evidence:** model mismatch is detectable at run-time

`http://probnum.org`          `https://pn.is.tue.mpg.de`

Probabilistic Numerics – Uncertainty in Computation
Hennig, Osborne, Girolami     Cambridge University Press, ETA 2019