



UPPSALA
UNIVERSITET

On the Construction of Probabilistic Newton-Type Algorithms

Thomas Schön, Uppsala University

Joint work with **Adrian Wills** at the University of Newcastle, Australia.

SIAM Conference on Uncertainty Quantification, Garden Grove, CA, USA.
April 18, 2018.

What? Solve the non-convex stochastic optimization problem

$$\max_x f(x)$$

when we only have access to **noisy** evaluations of $f(x)$ and its derivatives.

Why? These stochastic optimization problems are common:

- When the cost function cannot be evaluated on the entire dataset.
- When numerical methods approximate $f(x)$ and $\nabla^i f(x)$.
- ...

How? Learn a probabilistic nonlinear model of the Hessian.

Provides a local approximation of the cost function $f(x)$.

Use this local model to compute a search direction.

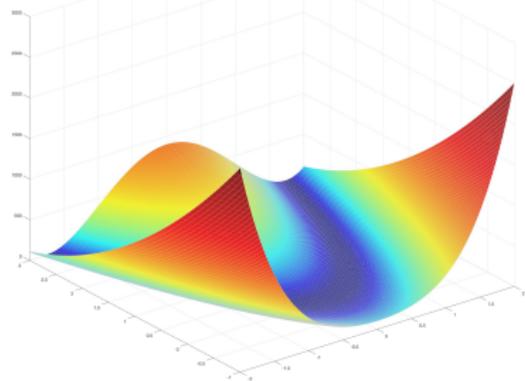
Captures second-order information (curvature) which opens up for better performance compared to a pure gradient-based method.

Intuitive preview example — Rosenbrock function

Let $f(x) = (a - x_1)^2 + b(x_2 - x_1^2)^2$, where $a = 1$ and $b = 100$.

Deterministic problem

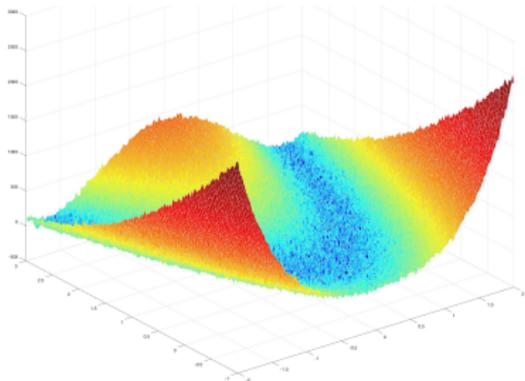
$$\min_x f(x)$$



Stochastic problem

$$\min_x f(x)$$

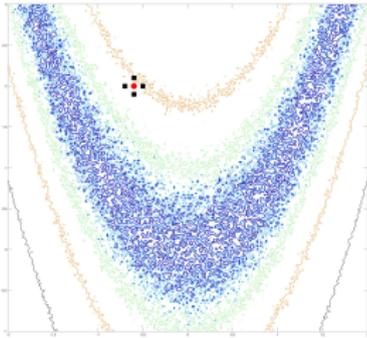
when we only have access to noisy versions of the cost function
 $(\tilde{f}(x) = f(x) + e, e \sim \mathcal{N}(0, 30^2))$
and its gradients.



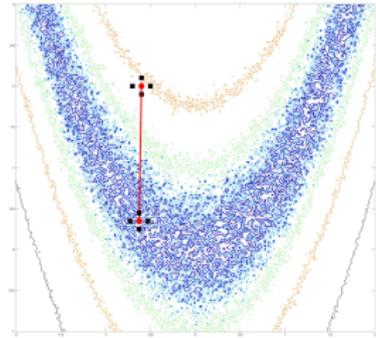
By not using the curvature information we expose ourself to the "banana-problem".

New algorithm at work — overall result

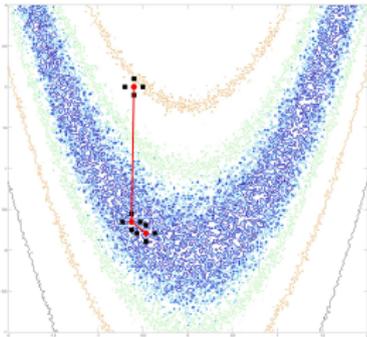
Initial value



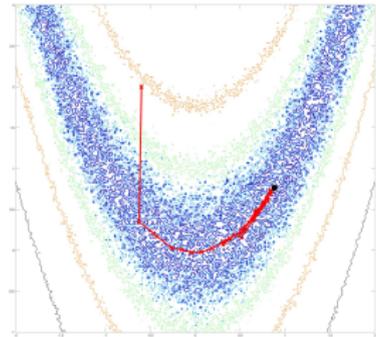
Iteration 1



Iteration 2



Iteration 50



Quasi-Newton — A non-standard take

Our problem is of the form

$$\max_{\mathbf{x}} f(\mathbf{x})$$

Idea underlying (quasi-)Newton methods: Learn a local quadratic model $q(\mathbf{x}_k, \delta)$ of the cost function $f(\mathbf{x})$ around the current iterate \mathbf{x}_k

$$q(\mathbf{x}_k, \delta) = f(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k)^\top \delta + \frac{1}{2} \delta^\top \mathbf{H}(\mathbf{x}_k) \delta$$

A second-order Taylor expansion around \mathbf{x}_k , where

$$\mathbf{g}(\mathbf{x}_k) = \nabla f(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{x}_k},$$

$$\mathbf{H}(\mathbf{x}_k) = \nabla^2 f(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{x}_k},$$

$$\delta = \mathbf{x} - \mathbf{x}_k.$$

We have measurements of the

- cost function $f_k = f(\mathbf{x}_k)$,
- and its gradient $\mathbf{g}_k = \mathbf{g}(\mathbf{x}_k)$.

Question: How do we update the Hessian model?

Line segment connecting two adjacent iterates \mathbf{x}_k and \mathbf{x}_{k+1} :

$$\mathbf{r}_k(\tau) = \mathbf{x}_k + \tau(\mathbf{x}_{k+1} - \mathbf{x}_k), \quad \tau \in [0, 1].$$

Useful basic facts

The fundamental theorem of calculus states that

$$\int_0^1 \frac{\partial}{\partial \tau} \nabla f(r_k(\tau)) d\tau = \nabla f(r_k(1)) - \nabla f(r_k(0)) = \underbrace{\nabla f(x_{k+1})}_{g_{k+1}} - \underbrace{\nabla f(x_k)}_{g_k}$$

and the chain rule tells us that

$$\frac{\partial}{\partial \tau} \nabla f(r_k(\tau)) = \nabla^2 f(r_k(\tau)) \frac{\partial r_k(\tau)}{\partial \tau} = \nabla^2 f(r_k(\tau)) (x_{k+1} - x_k).$$

$$\underbrace{g_{k+1} - g_k}_{=y_k} = \int_0^1 \frac{\partial}{\partial \tau} \nabla f(r_k(\tau)) d\tau = \int_0^1 \nabla^2 f(r_k(\tau)) d\tau \underbrace{(x_{k+1} - x_k)}_{s_k}.$$

Result — the quasi-Newton integral

With the definitions $y_k \triangleq g_{k+1} - g_k$ and $s_k \triangleq x_{k+1} - x_k$ we have

$$y_k = \int_0^1 \nabla^2 f(r_k(\tau)) d\tau s_k.$$

Interpretation: The difference between two consecutive gradients (y_k) constitute a *line integral observation of the Hessian*.

Problem: Since the Hessian is unknown there is no functional form available for it.

Solution 1 — recovering existing quasi-Newton algorithms

Existing quasi-Newton algorithms (e.g. BFGS, DFP, Broyden's method) assume the Hessian to be constant

$$\nabla^2 f(r_k(\tau)) \approx H_{k+1}, \quad \tau \in [0, 1],$$

implying the following approximation of the integral (**secant condition**)

$$y_k = H_{k+1} s_k.$$

Find H_{k+1} by **regularizing** H :

$$\begin{aligned} H_{k+1} &= \min_H \|H - H_k\|_W^2, \\ \text{s.t. } & H = H^\top, \quad H s_k = y_k, \end{aligned}$$

Equivalently, the existing quasi-Newton methods can be interpreted as **particular instances of Bayesian linear regression**.

Solution 2 — use a flexible nonlinear model

Our approach is fundamentally different.

Recall that the problem is **stochastic** and **nonlinear**.

Hence, we need a model that can deal with such a problem.

Idea: Represent the Hessian using a **Gaussian process** learnt from data.

Two of the remaining challenges:

1. Can we use line integral observations when learning a GP?
2. How do we ensure that the resulting GP represents a Hessian?

Stochastic quasi-Newton integral

$$y_k = \int_0^1 \underbrace{B(r_k(\tau))}_{=\nabla^2 f(r_k(\tau))} s_k d\tau + e_k,$$

corresponds to noisy (e_k) gradient observations.

Let us use a GP model for the unique elements of the Hessian

$$\tilde{B}(x) \sim \mathcal{GP}(\mu(x), \kappa(x, x')).$$

Resulting stochastic optimization algorithm

Standard non-convex numerical optimization loop with **non-standard components**.

Algorithm 1 Stochastic optimization

1. **Initialization** ($k = 1$)
 2. **while** *not terminated* **do**
 - (a) Compute a search direction p_k using the current approximation of the gradient g_k and Hessian B_k .
 - (b) Stochastic line search to find a step length α_k and set
$$x_{k+1} = x_k + \alpha_k p_k.$$
 - (c) Set $k := k + 1$
 - (d) Update the Hessian estimate (tailored GP regression)
 3. **end while**
-

Maximum likelihood nonlinear system identification

$$x_t = f(x_{t-1}, \theta) + w_t,$$

$$y_t = g(x_t, \theta) + e_t,$$

$$x_0 \sim p(x_0 | \theta),$$

$$(\theta \sim p(\theta)).$$

$$x_t | (x_{t-1}, \theta) \sim p(x_t | x_{t-1}, \theta),$$

$$y_t | (x_t, \theta) \sim p(y_t | x_t, \theta),$$

$$x_0 \sim p(x_0 | \theta),$$

$$(\theta \sim p(\theta)).$$

Maximum likelihood – model the unknown parameters as a deterministic variable θ and solve

$$\max_{\theta} p(y_{1:T} | \theta),$$

Challenge: The optimization problem is stochastic!

Cost function – the likelihood

Each element $p(y_t | y_{1:t-1}, \theta)$ in the likelihood

$$p(y_{1:T} | \theta) = \prod_{t=1}^T p(y_t | y_{1:t-1}, \theta),$$

can be computed by averaging over all possible values for the state x_t ,

$$p(y_t | y_{1:t-1}, \theta) = \int p(y_t | x_t, \theta) \underbrace{p(x_t | y_{1:t-1}, \theta)}_{\text{approx. by PF}} dx_t.$$

Non-trivial fact: The likelihood estimates obtained from the particle filter (PF) are **unbiased**.

Tutorial paper on the use of the PF (an instance of sequential Monte Carlo, SMC) for nonlinear system identification

ex) Simple linear toy problem

Identify the parameters $\theta = (a, c, q, r)^\top$ in

$$x_{t+1} = ax_t + w_t,$$

$$y_t = cx_t + e_t,$$

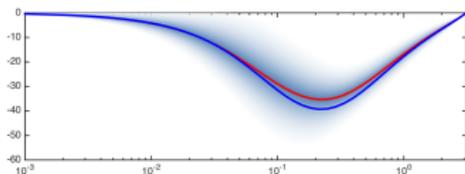
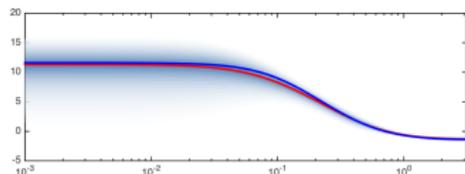
$$w_t \sim \mathcal{N}(0, q),$$

$$e_t \sim \mathcal{N}(0, r).$$

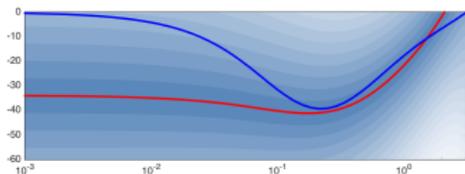
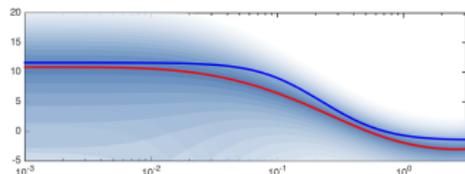
Observations:

- The likelihood $L(\theta) = p(y_{1:T} | \theta)$ and its gradient $\nabla_{\theta} L(\theta)$ are available in closed form via standard Kalman filter equations.
- Standard gradient-based search algorithms applies.
- Deterministic optimization problem $(L(\theta), \nabla_{\theta} L(\theta))$ noise-free).

ex) Simple linear toy problem



Both alg. in the noise-free case.



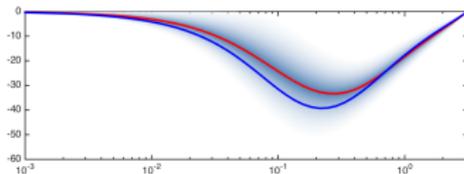
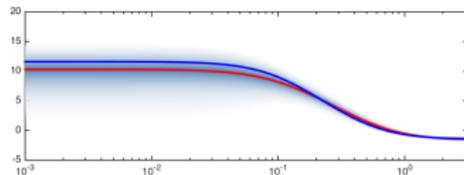
Classical BFGS alg. for noisy observations of $L(\theta)$ and $\nabla L(\theta)$.

100 independent datasets.

Clear blue – True system

Red – Mean value of estimate

Shaded blue – individual results



GP-based BFGS alg. with noisy observations of $L(\theta)$ and $\nabla L(\theta)$. 17/20

Ongoing work – scaling up to large problems

What is the key limitation of our GP-based optimization algorithm?

It **does not** scale to large-scale problems!

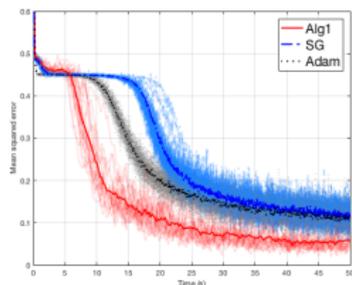
It is still highly useful and competitive for **small to medium** sized problems involving up to a coupled of hundred parameters or so.

We have developed a **new** technique that scales to **very large** problems.

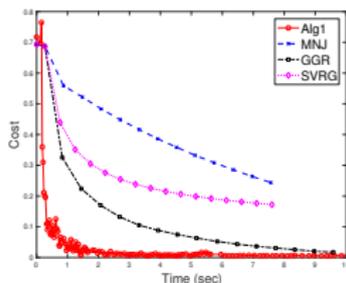
Ongoing work – scaling up to large problems

Key innovations:

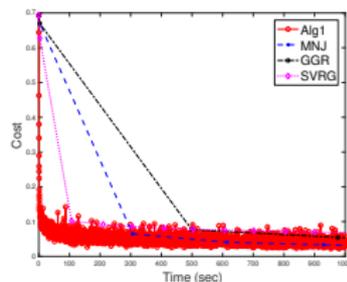
- Replace the GP with a matrix updated using fast Cholesky routines.
- Exploit a receding history of iterates and gradients akin to L-BFGS.
- An auxiliary variable Markov chain construction.



Training a deep CNN for MNIST data.



Logistic loss function with an L2 regularizer, gisette, 6 000 observations and 5 000 unknown variables.



Logistic loss function with an L2 regularizer, URL, 2 396 130 observations and 3 231 961 unknown variables.

Derived a **probabilistic** quasi-Newton algorithm that can be used with **noisy** observations of the cost function and its derivatives.

- Non-standard interpretation of quasi-Newton.
- Represent the Hessian using a Gaussian process.
- Application: Maximum likelihood estimation in nonlinear SSMs.
- We can scale up to large problems.

Adrian G. Wills and Thomas B. Schön. **On the construction of probabilistic Newton-type algorithms**, *Proceedings of the 56th IEEE Conference on Decision and Control (CDC)*, Melbourne, Australia, December 2017.

Significantly updated material will soon be available.