Gaussian Processes in Numerical Optimization

Philipp Hennig

GPSS, Sheffield 12 June 2013



Max Planck Institute for Intelligent Systems Department for Empirical Inference Research Group Elementary Intelligence

The machine learning pipeline

we need a coherent framework for hierarchical machine learning



- black boxes cause unaccounted errors
- making errors explicit may improve quality / lower cost
- joint language required: probability

Numerical algorithms are the elements of inference

inferring solutions of non-analytic problems

Numerical algorithms estimate an intractable property of a function given evaluations of function values.

quadrature estimate $\int_a^b f(x) dx$ given $\{f(x_i)\}$ analysis estimate x(t) under x' = f(x,t)given $\{f(x_i,t_i)\}$ control estimate $\min_u x(t,u)$ under x' = f(x,t,u) $\{f(x_i,t_i,u_i)\}$ optimization estimate $\arg \min_x f(x)$ given $\{f(x_i), \nabla f(x_i)\}$

- an inference problem
- even deterministic problems can be uncertain
- not a new idea¹, but rarely studied

¹H. Poincaré, 1896, Diaconis 1988, O'Hagan 1992

Numerical algorithms are learning machines!

now, let's focus on optimization

For $f : \mathbb{R}^N \to \mathbb{R}$, find local minimum $\arg \min f(x)$, starting at x_0 .

An old idea: Newton's method

$$f(x) \approx f(x_t) + (x - x_t)^{\mathsf{T}} \nabla f(x_t) + \frac{1}{2} (x - x_t)^{\mathsf{T}} \underbrace{\nabla \nabla^{\mathsf{T}} f(x_t)}_{=:B(x_t)} (x - x_t)$$
$$\underbrace{x_{t+1} = x_t - B^{-1}(x_t) \nabla f(x_t)}_{=:B(x_t)}$$

Cost: $\mathcal{O}(N^3)$

⊸⊳

High-dimensional optimization requires giving up knowledge in return for lower cost.

Quasi-Newton methods (think BFGS, DFP, ...)

aka. variable metric optimization - low rank estimators for Hessians

Instead of evaluating Hessian, build low-rank estimator

$$B_{t+1} = B_t + uv^{\mathsf{T}} + vu^{\mathsf{T}}$$
 with $u, v \in \mathbb{R}^{N \times M} \longrightarrow \mathcal{O}(N^2 + M^3)$

want estimator fulfilling local difference relation

$$\nabla f(x_{t+1}) - \nabla f(x_t) = B_{t+1}(x_{t+1} - x_t)$$
$$y_t = B_{t+1}s_t$$

- ... otherwise close to previous estimator in $||B_{t+1} B_t||_{F,V}$
- ... so minimize regularised loss

$$B_{t+1} = \underset{B \in \mathbb{R}^{N \times N}}{\arg \min} \left\{ \underset{\beta \to 0}{\lim} \frac{1}{\beta} \|y_t - Bs_t\|_V^2 + \|B - B_t\|_{F,V}^2 \right\}$$
$$= B_t + \frac{(y_t - B_t s_t)s_t^{\mathsf{T}}V}{s_t^{\mathsf{T}}V s_t}$$
Broyden's (1965) method

Quasi-Newton methods: the probabilist's view

Gaussian regression on matrices

P.H. & M. Kiefel, ICML 2012, JMLR 2013

$$B_{t+1} = \lim_{\beta \to 0} \arg\max_{B} \{-\frac{1}{\beta} \| y_t - Bs_t \|_{V}^2 - \| B - B_t \|_{F,V}^2 \}$$

$$= \lim_{\beta \to 0} \arg\max_{B} \{\log \underbrace{p(y_t | B, s_t, \beta, V)}_{\text{likelihood}} + \log \underbrace{p(B | B_t, V)}_{\text{prior}} \}$$

$$= \lim_{\beta \to 0} \arg\max_{B} p(y_t | B, s_t, \beta, V) p(B | B_t, V)$$

$$= \lim_{\beta \to 0} \arg\max_{B} \mathcal{N}(y_t; Bs_t, \beta V) \mathcal{N}(\overrightarrow{B}; \overrightarrow{B}_t, V \otimes V)$$

$$= \arg\max_{B} \mathcal{N}\left[B; B_t + \frac{(y_t - B_t s_t)Vs_t^{\mathsf{T}}}{s_t^{\mathsf{T}}Vs_t}, V \otimes \left(V - \frac{Vss^{\mathsf{T}}V}{s^{\mathsf{T}}Vs}\right)\right]$$

Quasi-Newton methods perform local maximum a-posteriori Gaussian inference on the Hessian's elements.

Outline

QN-methods perform regression!

Whenever you see ℓ_2 arguments, think of GPs!

what prior assumptions allow low cost?

Encoding symmetry exactly is expensive

full probabilistic solution has quartic cost

P.H., ICML 2013

the naïve idea:

• "observe symmetry" using $\Delta \vec{B} = \frac{1}{2} (\vec{B} - \vec{B^{\top}})$

$$p(0|B,\Delta) = \delta(0 - \Delta \overrightarrow{B})$$

• however, $\operatorname{rk} \Delta = \frac{1}{2}N(N-1)$: costs at least $\mathcal{O}(N^4)$

Instead, make two independent observations

observe \overline{B} acting on two different spaces

instead: Powell's (1970) Symmetric Broyden (PSB) update

- before: y = Bs
- now, additionally: "dual observation"

$$p(y_i^{\mathsf{T}} | B, s_i^{\mathsf{T}}, V_{i-1}) = \delta(y_i^{\mathsf{T}} - s_i^{\mathsf{T}} B) = \lim_{\beta \to 0} \mathcal{N}(y_i^{\mathsf{T}}; s_i^{\mathsf{T}} B, \beta \otimes V_i)$$

gives posterior mean, covariance

$$B_{i} = B_{i-1} + \frac{(y_{i} - B_{i-1}s_{i})s_{i}^{\mathsf{T}}V_{i-1} + V_{i-1}s_{i}(y_{i} - B_{i-1}s_{i})^{\mathsf{T}}}{s_{i}^{\mathsf{T}}V_{i-1}s_{i}} - \frac{V_{i-1}s_{i}(s_{i}^{\mathsf{T}}(y_{i} - B_{i-1}s_{i}))s_{i}^{\mathsf{T}}V_{i-1}}{(s_{i}^{\mathsf{T}}V_{i-1}s_{i})^{2}}$$
$$\Sigma_{i} = \left(V_{i-1} - \frac{V_{i-1}s_{i}s_{i}^{\mathsf{T}}V_{i-1}}{s_{i}^{\mathsf{T}}V_{i-1}s_{i}}\right) \otimes \left(V_{i-1} - \frac{V_{i-1}s_{i}s_{i}^{\mathsf{T}}V_{i-1}}{s_{i}^{\mathsf{T}}V_{i-1}s_{i}}\right)$$

algebraic assumptions used to shape computational cost

$$p(B) = \mathcal{N}(B; B_0, V \otimes V)$$

$$\operatorname{cov}(B_{ij}, B_{k\ell}) = (V \otimes V)_{(ij), (k\ell)} = V_{ik}V_{j\ell}$$

$$= \operatorname{cov}(\partial_i(\nabla f), \partial_k(\nabla f)) \cdot \operatorname{cov}(\partial_j(f \nabla^{\mathsf{T}}), \partial_\ell(f \nabla^{\mathsf{T}}))$$

- each $B_{i:}$ is a map $\mathbb{R}^{N \times 1} \rightarrow \mathbb{R}$: $\partial_i (\nabla f)$
- each $B_{:j}$ is a map $\mathbb{R}^{1 \times N} \to \mathbb{R}$: $\partial_j (f \nabla^{\mathsf{T}})$
- dual vector spaces are isomorphic, thus $(\nabla f)\nabla^{\top} = \nabla (f\nabla^{\top})$
- but we won't tell the method! Treat the two maps as independent.
- replace $f(x) : \mathbb{R}^N \to \mathbb{R}$ with $f(x^{\mathsf{T}}, x) : \mathbb{R}^{N \times 1} \times \mathbb{R}^{1 \times N} \to \mathbb{R}$

Forgetting your good education to get ahead

QN methods deliberately ignore algebraic knowledge to lower computational cost



To achieve low computational cost, numerical methods may have to deliberately ignore prior knowledge!

Whenever you see Kronecker products, someone is hiding an independence assumption!

where are the GPs already?

Hessians are not constant

can we have a prior over regular Hessians?

The unavoidable slide

nonparametric generalisation of Gaussians to space of real-valued functions

$$p(f) = \mathcal{GP}(f; \mu, k)$$
$$p(Af) = \mathcal{GP}(Af; A\mu, AkA^{\mathsf{T}})$$

The Gaussian family is closed under linear projections

thus provides the functionality for quasi-Newton inference

$$p(f) = \mathcal{GP}(f; \mu, k)$$

$$p(Af) = \mathcal{GP}(Af; A\mu, AkA^{\mathsf{T}})$$

$$p\left(\int_{a}^{b} f(x)dx\right) = \mathcal{GP}\left(\int_{a}^{b} f(x)dx; \int_{a}^{b} \mu(x)dx, \iint_{a}^{b} k(x, x')dxdx'\right)$$

Nonparametric quasi-Newton methods

Inferring a changing Hessian

P.H. & M. Kiefel, ICML 2012, JMLR 2013

Idea: replace

$$\nabla f(x_{t+1}) - \nabla f(x_t) \approx B(x_{t+1} - x_t)$$
$$\Rightarrow = \int_{x_t}^{x_{t+1}} B(x) \, dx$$

• Gaussian process prior on $B(x^{\mathsf{T}}, x)$

$$p(B) = \mathcal{GP}(B, B_0(x^{\mathsf{T}}, x), k(x^{\mathsf{T}}, {x'}^{\mathsf{T}}) \otimes k(x, x'))$$

Gaussian likelihoods

$$p(y_i(x^{\mathsf{T}}) | B, s_i) = \lim_{\beta \to 0} \mathcal{N}\left(y_i; \sum_m s_{im} \int_0^1 B(x^{\mathsf{T}}, x(t)) \, \mathrm{d}t, k(x^{\mathsf{T}}, {x'}^{\mathsf{T}}) \otimes \beta\right)$$
$$p(y_i(x)^{\mathsf{T}} | B, s_i^{\mathsf{T}}) = \lim_{\beta \to 0} \mathcal{N}\left(y_i^{\mathsf{T}}; \sum_m s_{im}^{\mathsf{T}} \int_0^1 B(x^{\mathsf{T}}(t), x) \, \mathrm{d}t, \beta \otimes k(x, x')\right)$$

- posterior of same algebraic form as before, but with linear maps of nonlinear (integral of k) entries.
- same computational complexity as L-BFGS (Nocedal, 1980): $\mathcal{O}(N)$

nonparametric inference on elements of the Hessian



nonparametric inference on elements of the Hessian



nonparametric inference on elements of the Hessian



nonparametric inference on elements of the Hessian



nonparametric inference on elements of the Hessian



Advantages of nonparametric modelling

analytic use of all observations

- choice of kernel allows more flexible model
- exact treatment of gradient observations
- can use all observations from line searches



Some empirical results

optimizing a 200-dimensional function

P.H. & M. Kiefel, ICML 2012, JMLR 2013



23

quasi-Newton methods can be generalised to GP models.

GPs are cubic in the number of inputs, not the outputs!

how does this all help machine learning?

optimization of noisy objectives

- noisy objectives a frequent in learning (mini-batch training)
- stochastic gradient descent widely used
- noisy quasi-Newton methods improve conceptually



optimization of dynamically changing objectives

further generalisations

► dynamic settings → filtering



- the numerical methods used in ML are learners themselves don't treat them as black boxes
- ▶ BFGS is a Gaussian regressor think $\|\cdot\|_2 \leftrightarrow \mathcal{N}(x)$ and $\|\cdot\|_F \leftrightarrow \mathcal{N}(A)$ and $\|\cdot\|_{HS} \leftrightarrow \mathcal{GP}(f)$
- ► to be fast, you may need to ignore even pertinent knowledge watch out for ⊗ in covariances
- GP numerical optimization need not be expensive

design your model carefully

- machine learning needs probabilistic numerical methods
- major theoretical, practical questions still open

Want to stay informed? Befriend Maren!

Are quasi-Newton methods locally linearly convergent?

Definition (local linear convergence)

QN update is locally linearly convergent at the minimum x_* if $\exists \epsilon, \delta > 0$, s.t.

$$(\|x_0 - x_*\| < \epsilon \land \|B_0 - B(x_*)\| < \delta) \implies (\|x_k - x_*\| = \mathcal{O}(e^{-k}))$$

nonparametric QN methods does not, in general, obey this, just like nonlinear regression does not collapse linearly fast.

Theorem (Broyden, Dennis, Moré, 1973)

QN update is locally linearly convergent if

$$||B_i - B(x_*)|| \le [1 + \alpha \sigma(x_{i-1}, x_i)] ||B_{i-1} - B(x_*)|| + \beta \sigma(x_{i-1}, x_i)$$

$$\sigma(x_{i-1}, x_i) \coloneqq \max\{||x_i - x_*||, ||x_{i-1} - x_*||\}$$

positive definiteness?

as similar story to symmetry

the naïve idea:

prior exclusively over positive definite matrices: Wishart

 $B = uu^{\mathsf{T}}$ with $u_i \sim \mathcal{N}(0, V)$

however, inference not analytically tractable!

as similar story to symmetry

instead: BFGS / DFP:

$$p(B) = \mathcal{N}(B; B_0, V \otimes V) \quad \text{and set} \quad V = B$$
$$\propto |B|^{-N^2/2} \cdot \exp\left[-\frac{1}{2}\left(N - 2\operatorname{tr}(B_0 B^{-1}) + \operatorname{tr}(B_0 B^{-1} B_0 B^{-1})\right)\right]$$

- posterior mean is positive definite if $y^{\mathsf{T}}s > 0$
- but prior puts nonzero mass on indefinite matrices like

$$B = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \qquad \Rightarrow \qquad p(B) \propto \exp(-2)$$

BFGS and DFP have other advantages (scale-freedom, etc.). But they do not explicitly encode positive definiteness.

of optimization steps, 10^3 samples



of optimization steps, $3 \cdot 10^4$ samples



of optimization steps, $6 \cdot 10^4$ samples



wallclock time, 10^3 samples



wallclock time, $3\cdot 10^4~\text{samples}$



wallclock time, $6\cdot 10^4$ samples

