

Optimizing Magnetic Confinement Devices for Fusion Plasmas

David Bindel

12 Sep 2022

Department of Computer Science
Cornell University

Who?

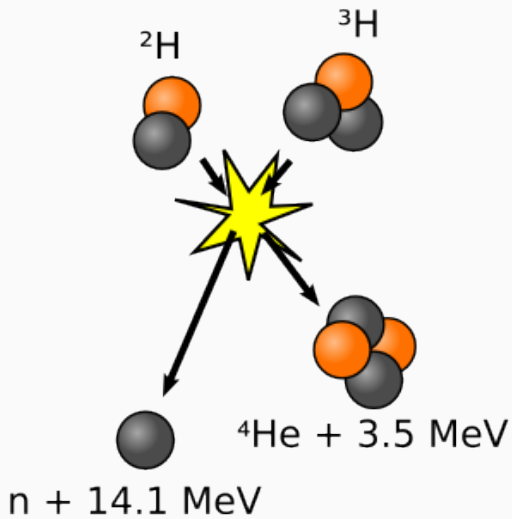
Simons Collaboration: “Hidden Symmetries and Fusion Energy”

<https://hiddensymmetries.princeton.edu/>

Princeton, NYU, Maryland, IPP Greifswald, Warwick, CU Boulder, UW Madison, EPFL, ANU, UT Austin, U Arizona.

- Phase 1: Sep 2017-Aug 2022
- Phase 2: Sep 2022-Aug 2025

D-T fusion



Lawson criterion

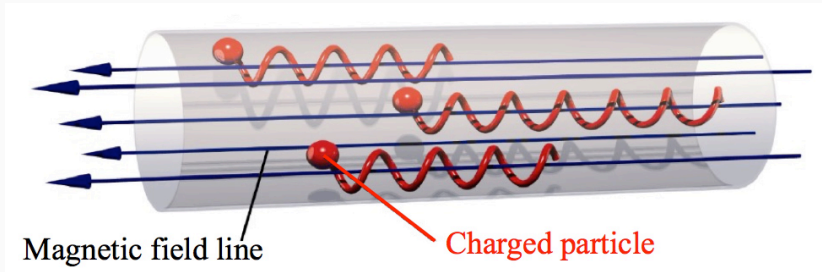
Figure of merit: $nT\tau_E$ where

- n is number density
- T is temperature
- τ_E is energy confinement time

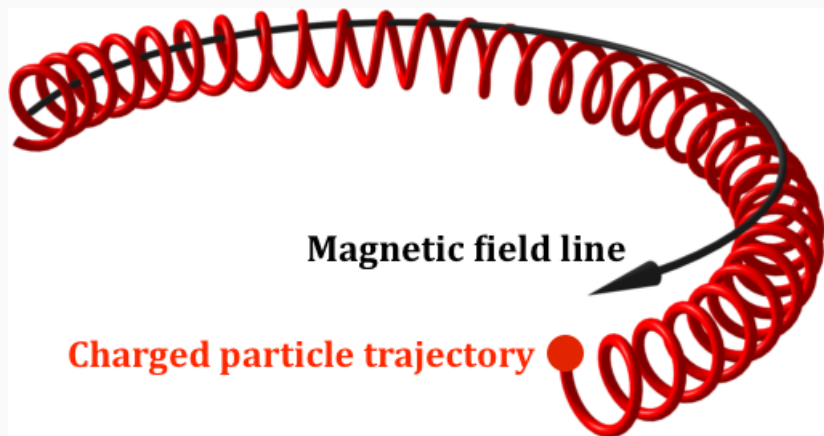
Min value required at $T = 14$ keV (about 162×10^6 K) is

$$nT\tau_E \geq 3.5 \times 10^{28} \text{ K s/m}^3$$

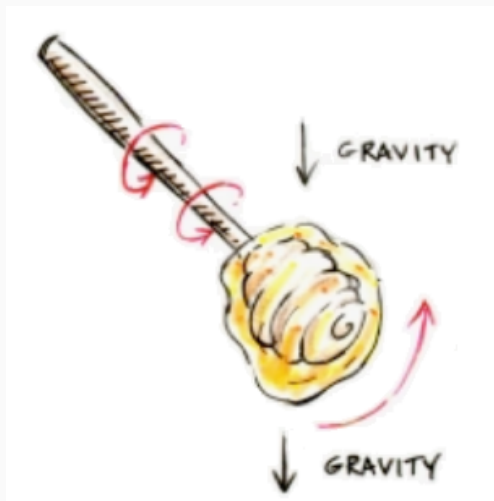
Magnetic confinement basics



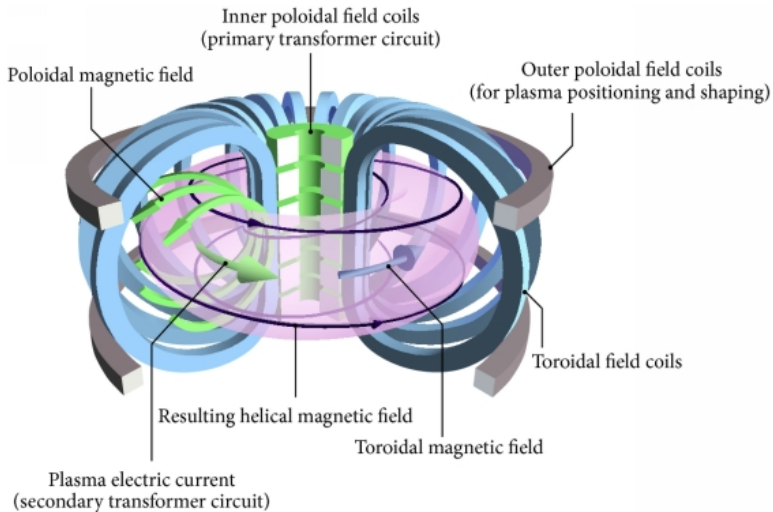
Magnetic confinement basics

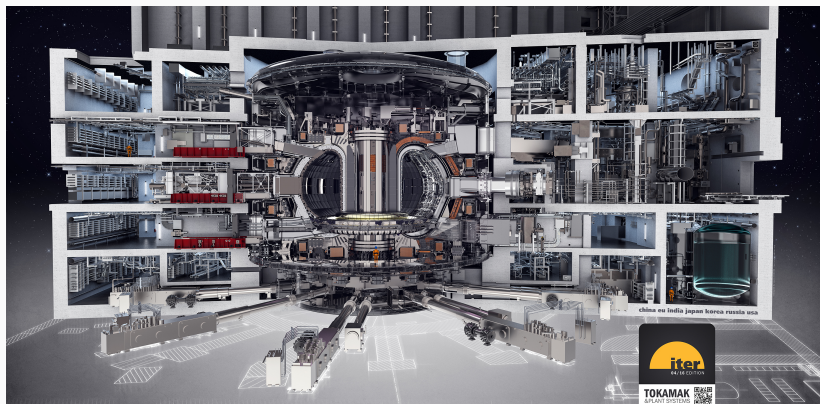


Magnetic confinement basics



The big name: Tokamaks

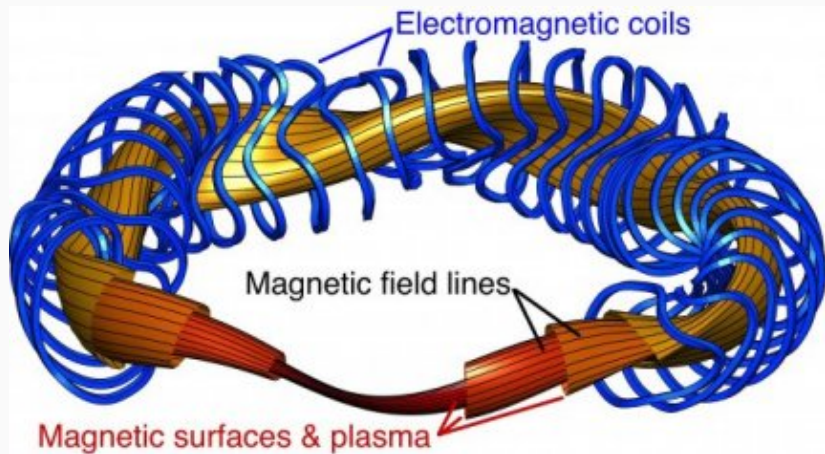




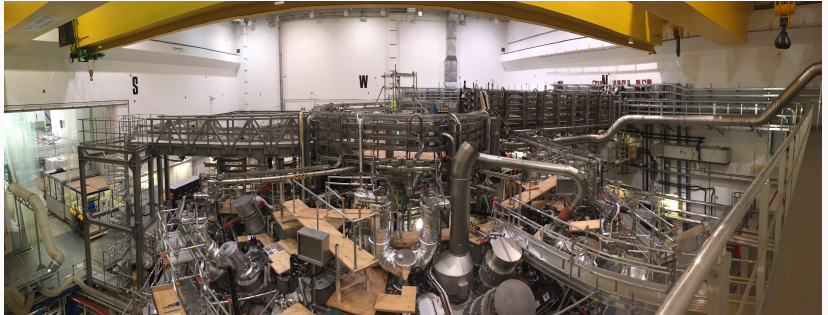
Trouble with Tokamaks

- Toroidal and poloidal coils produce a helical B
- Also a strong current (which affects the field)
- Prone to *disruption* (a bad instability)

Stellarator Concept



Wendelstein 7-X Machine



Operating since 2015-12-10;
plasma discharges lasting several min.

Trouble with Stellarators

- Stability is not such an issue
- More “neo-classical” (low-collisionality) transport — bad for confinement
- Control neo-classical transport by quasisymmetry ($|B|$ has a symmetry in Boozer coordinates)

Stellarator Quality Measures

What makes an “optimal” stellarator?

- Approximates field symmetries (which measures?)
- Satisfies macroscopic and local stability
- Divertor fields for particle and heat exhaust
- Minimizes collisional and energetic particle transport
- Minimizes turbulent transport
- Satisfies basic engineering constraints (cost, size, etc)

Each objective involves different approximations, uncertainties, and computational costs.

What Makes a Good Stellarator?

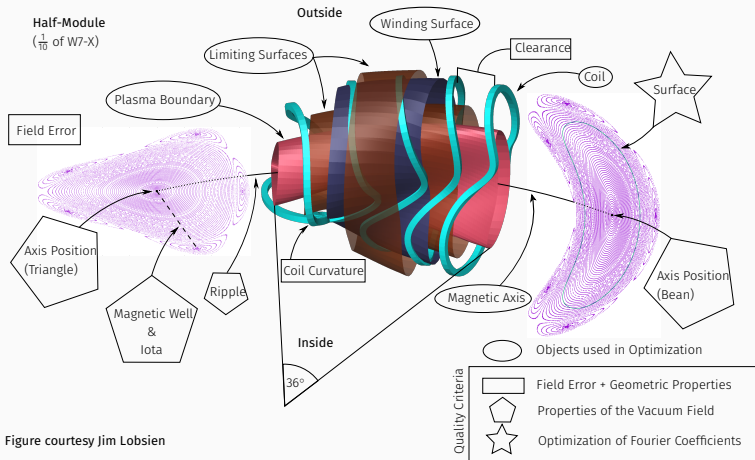
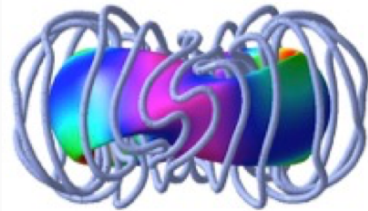
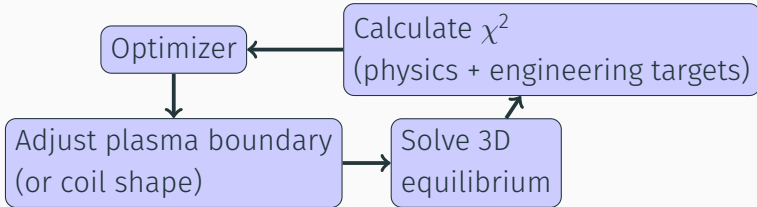
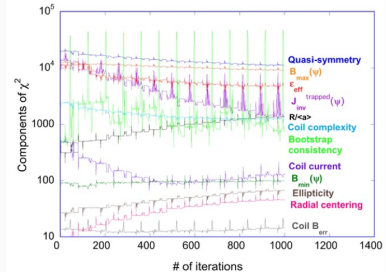


Figure courtesy Jim Lobsien

How Do We Optimize? (STELLOPT Approach)



$$r(\phi, \theta) + iz(\phi, \theta) = \sum \alpha_{m,n} e^{i(m\phi - n\theta)}$$



Challenges

1. Costly and “black box” physics computations
 - Each step: MHD equilibrium solve, transport, coil design, ...
 - Several times per step for finite-difference gradients
2. Managing tradeoffs
 - How do we choose the weights in the χ^2 measure? By gut?
 - Varying the weights does not expose tradeoffs sensibly
3. Dealing with uncertainties
 - What you simulate \neq what you build!
4. Global search
 - How to avoid getting stuck in local minima?

Progress of the Simons collaboration

- Collaboration has made a lot of progress (though work remains) on
 - Fast equilibrium solvers (NYU, Arizona, Flatiron)
 - Faster simulations, with derivatives (NYU, Maryland, Princeton)
 - Optimizing under uncertainty (Greifswald, Cornell)
 - Producing plasmas with high quasisymmetry (Maryland, Princeton)
- More limited progress on
 - Global search (though some with TuRBO)
 - Fast and accurate proxies for turbulent transport
 - Optimizing with instabilities (micro/macro)
 - Optimization of divertors

Rest of the session on two Cornell-centered projects:

- Optimization under uncertainty
- Multi-objective optimization

Optimization Under Uncertainty

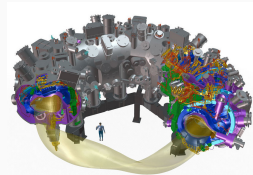
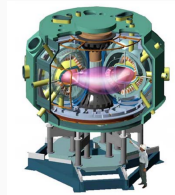
Low construction tolerances:

- NCSX: 0.08%
- Wendelstein 7-X: 0.1% – 0.17%

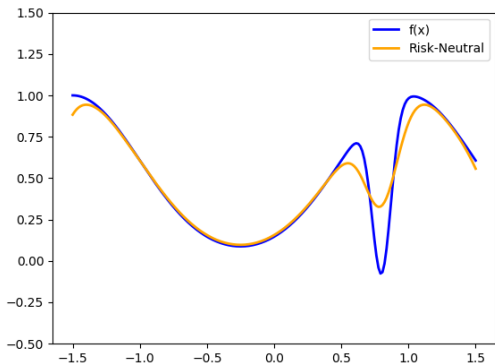
Higher tolerances as coil opt goal!

Also want tolerance to

- Changes to control parameters
- Uncertainty in physics or model



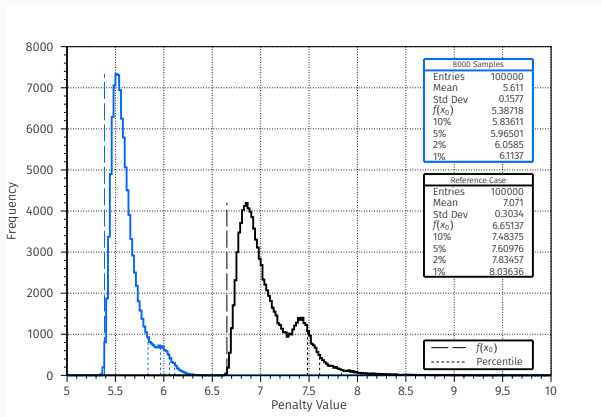
Risk-neutral OUU



Want efficient OUU in ~ 200 dimensions

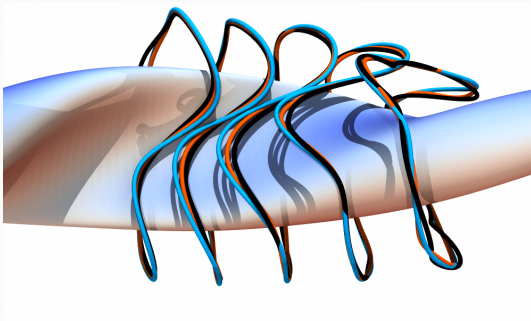
$$\min_{x \in \Omega} \mathbb{E}_U[f(x - U)]$$

(Recent) Prior: Monte Carlo Approach



Robustness & mean perf greatly improved (w/ $\sim 10^8$ evals)
J.-F. Lobsien, M. Drevlak, T. Kruger, S. Lazerson, C. Zhu, T. S. Pedersen,
Improved performance of stellarator coil design optimization,
Journal of Plasma Physics, 2020.

Our Approach: fast TuRBO-ADAM



Black: ref; red: TuRBO-ADAM 10mm; blue: TuRBO-ADAM 20mm.

Evaluate objective with FOCUS from PPPL.

- Global search with modified TuRBO
- Local refinement with ADAM with control variate

Costs about 0.01% the evaluation budget.

Interpolation / regression / supervised learning:

- Choose an approximation family \mathcal{F}
 - Global bases (Fourier, polynomial, etc)
 - Local bases (e.g. piecewise polynomials)
 - Gaussian processes or splines
 - Neural networks
 - etc
- Choose $s \in \mathcal{F}$ to agree with data about f
 - May need regularization / priors
 - May acquire data adaptively
- Use s in lieu of f

Error bounds: Consistency (how close f is to \mathcal{F}) + stability

Curse of Dimensionality

Suppose $f : \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}$ has m derivatives. Worst-case approximation error over Ω scales with measurements n like

$$\|f - \hat{f}\|_{\infty} \sim \text{diam}(\Omega)n^{-m/d}.$$

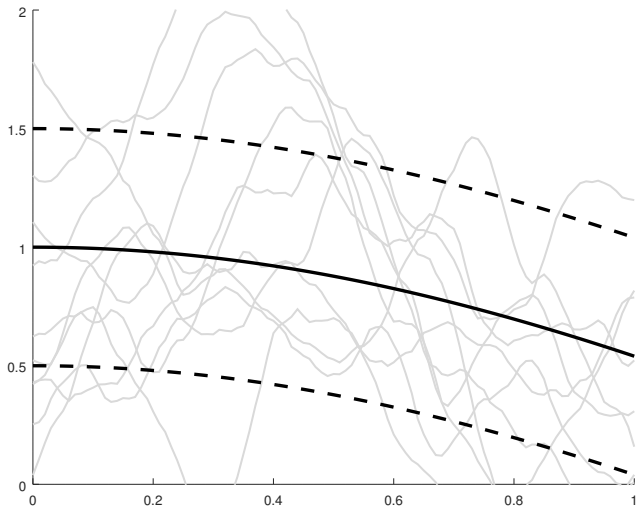
Optimal rate for *global* opt differs by a constant.
(Local optimization is a different problem.)

Sample efficiency requires:

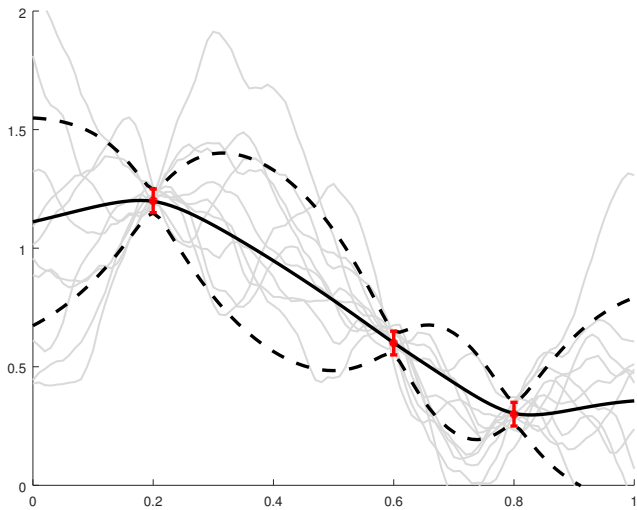
- Lots of smoothness (large m)
- Low-dimensional structure (small “effective” d)
- Or some other structure (sometimes w/o justification)

Usually *do* have more structure...

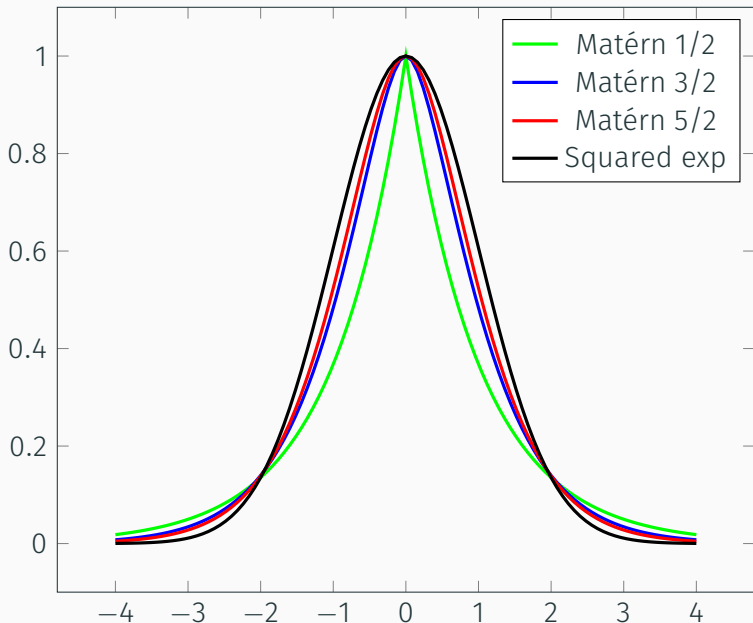
Gaussian Processes (GPs)



Being Bayesian



Matérn and SE kernels



Gaussian Processes (GPs)

Our favorite continuous distributions over

$$\mathbb{R}: \quad \text{Normal}(\mu, \sigma^2), \quad \mu, \sigma^2 \in \mathbb{R}$$

$$\mathbb{R}^n: \quad \text{Normal}(\mu, C), \quad \mu \in \mathbb{R}^n, C \in \mathbb{R}^{n \times n}$$

$$\mathbb{R}^d \rightarrow \mathbb{R}: \quad \text{GP}(\mu, k), \quad \mu : \mathbb{R}^d \rightarrow \mathbb{R}, k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

More technically, define GPs by looking at finite sets of points:

$$\forall X = (x_1, \dots, x_n), x_i \in \mathbb{R}^d,$$

have $f_X \sim N(\mu_X, K_{XX})$, where

$$f_X \in \mathbb{R}^n, \quad (f_X)_i \equiv f(x_i)$$

$$\mu_X \in \mathbb{R}^n, \quad (\mu_X)_i \equiv \mu(x_i)$$

$$K_{XX} \in \mathbb{R}^{n \times n}, \quad (K_{XX})_{ij} \equiv k(x_i, x_j)$$

When X is unambiguous, we will sometimes just write K .

Being Bayesian

Now consider prior of $f \sim \text{GP}(\mu, k)$, noisy measurements

$$f_X \sim y + \epsilon, \quad \epsilon \sim N(0, W), \quad \text{typically } W = \sigma^2 I$$

Posterior is $f \sim \text{GP}(\mu', k')$ with

$$\begin{aligned} \mu'(x) &= \mu(x) + K_{xx}c & \tilde{K} &= K_{xx} + W \\ k'(x, x') &= K_{xx'} - K_{xx}\tilde{K}^{-1}K_{xx'} & c &= \tilde{K}^{-1}(y - \mu_X) \end{aligned}$$

The expensive bit: solves with \tilde{K} .

Bayesian Optimization (BO)

Typical GP-based BO:

- Evaluate f on initial sample in Ω
- Condition a GP on sample data
- Until budget exhausted
 - Optimize *acquisition function* $\alpha(x)$ over Ω
(e.g. $\alpha_{\text{EI}}(x) = E [[f(x_{\text{best}}) - f(x)]_+]$ where x_{best} is best so far)
 - Evaluate at selected point
 - Update the GP model (including hyper-parameters)
- Standard cost: $O(n^3)$ per step (with n data points)

Great for modest budgets — but no escape from the Curse!

To avoid the Curse, we need an assumption!

So: suppose d large, but not too many minimizers:

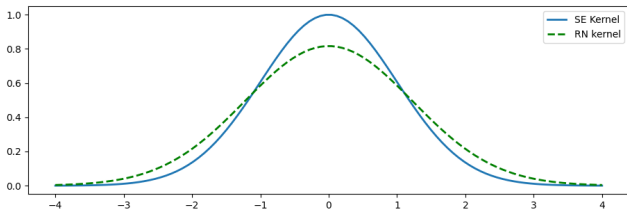
- Choose M starting points scattered over Ω
- Run local minimizer (gradient descent, Newton, etc)
- Hope for at least one start per convergence basin

Q: How to allocate effort to different starts?

For high-d: combine local BO with multi-start strategy

- Rough global sampling at M points
- Local GP models and trust-region around each point
- Thompson sampling to choose which local model (and trust region) to refine next

(Eriksson, Pearce, Gardner, Turner, Poloczek, 2019)



- TuRBO builds GP models for $f(x)$ (nominal objective)
- Simple transform from GP for $f(x)$ to GP for $E_U[f(x + U)]$ (Beland and Nair, 2017)

Problem: TuRBO explores a lot — want more refinement

Stochastic Gradient Descent (SGD)

Ordinary gradient descent is

$$x_{k+1} = x_k - \alpha_k \nabla \phi(x_k)$$

SGD is

$$x_{k+1} = x_k - \alpha_k g_k$$

where g_k is a random draw, $E[g_k] = \nabla \phi(x_k)$.

For $\phi(x) = E_U[f(x + U)]$, use $g_k = \nabla f(x_k + u_k)$.

Convergence is slow ($O(1/m)$), but steps can be cheap.
Speed depends a lot on variance of gradient estimator.

Adam + Control Variates

- Regular Adam: stochastic gradient algorithm with “adaptive momentum” for step size control. Use directions

$$g(x) = \nabla f(x + U)$$

for a random draw U (can also do mini-batch).

- Variance reduction with control variates (Wang, Chen, Smola, Xing, 2013)

$$g(x) = \nabla f(x + U) + \alpha(\hat{g}(x) - E[\hat{g}(x)])$$

$$\hat{g}(x) = \nabla f(x) + HU.$$

- True Hessian not avail, so set H to be an approximate Hessian (BFGS approximation via gradients from Adam).

Additional Information

Multi-output GPs model $f: \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}^k$

- Model covariance over space and across outputs.
- Example: function values + derivatives

$$\mu^\nabla(\mathbf{x}) = \begin{bmatrix} \mu(x) \\ \nabla_x \mu(x) \end{bmatrix}, \quad k^\nabla(x, x') = \begin{bmatrix} k(x, x') & (\nabla_{x'} k(x, x'))^T \\ \nabla_x k(x, x') & \nabla^2 k(x, x') \end{bmatrix}$$

- Can also model multi-fidelity sims, etc

Pro: FOCUS provides gradients, easy to incorporate!

Con: Matrix dimensions scale like $n(d + 1)$

Idea: *variational* approximate inference with derivatives

- Assume posterior of a given functional form
- Minimize evidence lower bound (ELBO) via Adam
- Demo with half a million FOCUS runs ($n = 500K$, $d = 45$)

(Bindel, Gardner, Huang, Padidar, Zhu, NeurIPS 2021)

Constrained and Multi-Objective

Naive: put everything we care about in a nonlinear LS problem

- $f_k(x)$ is deviation from k th target
- Add some weighting (chosen by the user)

But is this actually what we want?

- Choice of target values is unclear
- Choice of weights is unclear

And there are reasons for numerical nervousness:

- Maybe too few objectives (underdetermined LS problems)
- Maybe poorly conditioned (esp. with “large” weights)
- May not have small residual

Tackling Constraints

General problem

$$\text{minimize } \phi(x) \text{ s.t. } \begin{cases} c_j(x) = 0, & j \in \mathcal{E} \\ c_j(x) \leq 0, & j \in \mathcal{I} \end{cases}$$

Convert into unconstrained optimization / nonlinear equation solving problem with:

- Fewer degrees of freedom (constraint elimination)
- Same degrees of freedom (penalties and barriers)
- More degrees of freedom (Lagrange multipliers)

Constraint elimination usually only for linear constraints.

KKT Conditions

$$\text{minimize } \phi(x) \text{ s.t. } \begin{cases} c_j(x) = 0, & j \in \mathcal{E} \\ c_j(x) \leq 0, & j \in \mathcal{I} \end{cases}$$

Define the Lagrangian

$$L(x, \lambda, \mu) = \phi(x) + \sum_{i \in \mathcal{E}} \lambda_i c_i(x) + \sum_{i \in \mathcal{I}} \mu_i c_i(x).$$

KKT conditions are

$$\nabla_x L(x^*) = 0$$

$$c_i(x^*) = 0, \quad i \in \mathcal{E} \quad \text{equality constraints}$$

$$c_i(x^*) \leq 0, \quad i \in \mathcal{I} \quad \text{inequality constraints}$$

$$\mu_i \geq 0, \quad i \in \mathcal{I} \quad \text{non-negativity of multipliers}$$

$$c_i(x^*) \mu_i = 0, \quad i \in \mathcal{I} \quad \text{complementary slackness}$$

Penalties and Barriers

Want to minimize

$$\text{minimize } \phi(x) \text{ s.t. } \begin{cases} c_j(x) = 0, & j \in \mathcal{E} \\ c_j(x) \leq 0, & j \in \mathcal{I} \end{cases}$$

Instead minimize for small γ

$$\psi_\gamma(x) = \phi(x) + \frac{1}{2\gamma} \sum_{i \in \mathcal{E}} c_i(x)^2 - \gamma \sum_{i \in \mathcal{I}} \log(-c_i(x)).$$

Note that at minimizer x^* :

$$\nabla \psi_\gamma(x^*) = \nabla \phi(x^*) + \sum_{i \in \mathcal{E}} \tilde{\lambda}_i \nabla c_i(x^*) + \sum_{i \in \mathcal{I}} \tilde{\mu}_i \nabla c_i(x^*)$$

where Lagrange multiplier estimates come from the c_i :

$$\tilde{\lambda}_i = c_i(x^*)/\gamma, \quad \tilde{\mu}_i = \gamma/c_i(x^*)$$

Standard trick: Penalty to estimate multipliers.

What about using nonlinear least squares for tradeoffs?

More generally, consider $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, maybe minimize

$$w^T f(x) = \sum_{k=1}^m w_k f_k(x).$$

Structural Optimization 14, 63–69 © Springer-Verlag 1997

A closer look at drawbacks of minimizing weighted sums of objectives for Pareto set generation in multicriteria optimization problems

I. Das and J.E. Dennis

Department of Computational and Applied Mathematics, Rice University of Houston, TX 77251-1892, USA

June 4, 2015

Matt Landreman

Some optimal solutions to a smooth multi-objective problem cannot be found by minimizing a total χ^2

Exploring the Pareto Frontier

x dominates y if

$$\forall k, f_k(x) \leq f_k(y)$$

and not all strict.

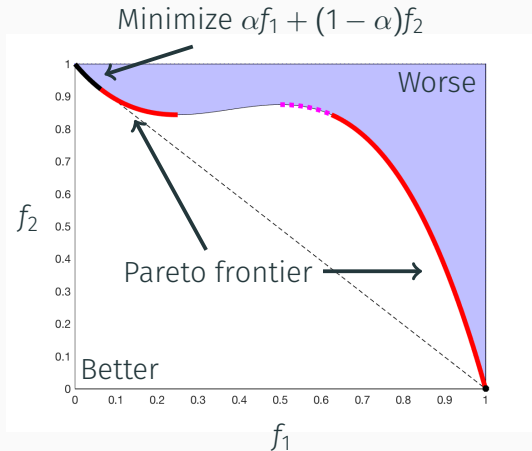
Best points are:

Pareto optimal,
aka non-dominated,
aka non-inferior,
aka non-efficient.

Form **Pareto frontier**

Minimizing $\sum_k \alpha_k f_k$ only explores convex hull!

Other methods sample / approximate the full frontier.



First-order condition

Stationary condition:

$$\{J(x)u : u \geq 0\} \cap \mathbb{R}_+^n = \{0\}.$$

Fritz John stationary condition: for some $\lambda \geq 0, \lambda \neq 0$

$$J(x)^T \lambda = 0.$$

Follows via Motzkin's theorem of the alternative: if A and C are given matrices, can either solve

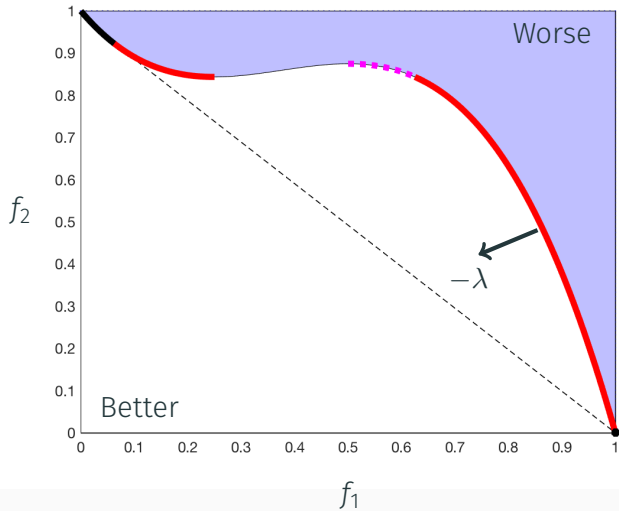
$$Ax < 0, \quad Cx \leq 0$$

or

$$A^T \lambda + C^T \mu = 0, \quad \lambda \geq 0, \lambda \neq 0, \mu \geq 0$$

But not both.

Fritz John multiplier geometry



Fritz John vs KKT

Fritz John condition (with constraints): Weak Pareto for

$$\text{minimize } f(x) \text{ s.t. } c(x) \leq 0$$

requires $\lambda \geq 0$ and $\mu \geq 0$ not both all zero such that

$$\lambda^T f'(x^*) + \mu^T c'(x^*) = 0$$

$$\mu_i c_i(x^*) = 0$$

Very similar to KKT conditions for constrained opt:

$$\nabla_x L(x^*) = 0, \quad L(x, \lambda, \mu) = \phi(x) + \lambda^T c_{\mathcal{E}}(x) + \mu^T c_{\mathcal{I}}(x)$$

$$c_i(x^*) = 0, \quad i \in \mathcal{E} \quad \text{equality constraints}$$

$$c_i(x^*) \leq 0, \quad i \in \mathcal{I} \quad \text{inequality constraints}$$

$$\mu_i \geq 0, \quad i \in \mathcal{I} \quad \text{non-negativity of multipliers}$$

$$c_i(x^*) \mu_i = 0, \quad i \in \mathcal{I} \quad \text{complementary slackness}$$

Constrained vs multi-objective

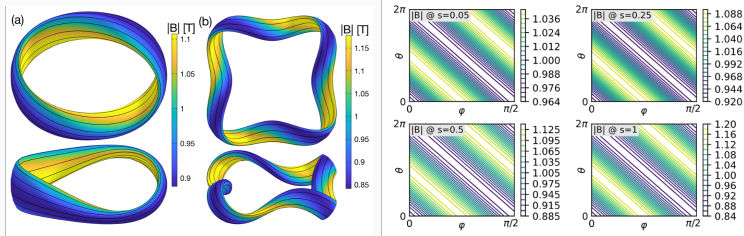
- First-order conditions are *almost* the same
- Can mix and match (constrained multi-objective)
- Multi-objective involves many solves to explore space
- Curse of dimensionality: exploration cost scales exponentially with m

Scalarizing

Find Pareto points via a single-objective optimization problem:

- Linear: $\phi(x) = w^T f(x)$
 - Need to consider stationary points to get full frontier.
 - Uniform weight sampling \neq uniform frontier sampling.
- Projection: $\phi(x) = \sum_i w_i (f_i(x) - f_i^*)^2$
 - Effectively what is done now.
 - Similar tradeoffs to linear scalarization.
- Chebyshev: $\phi(x) = \max_i w_i f_i(x)$
 - Nonsmooth where max is non-unique.
 - Uniform weight \neq uniform frontier sampling.
- ϵ -constraint: $\phi(x) = f_i(x), f_j(x) \leq \epsilon_j$ for $j \neq i$
 - Subproblem is constrained.
 - Can get uniform sampling in components other than i

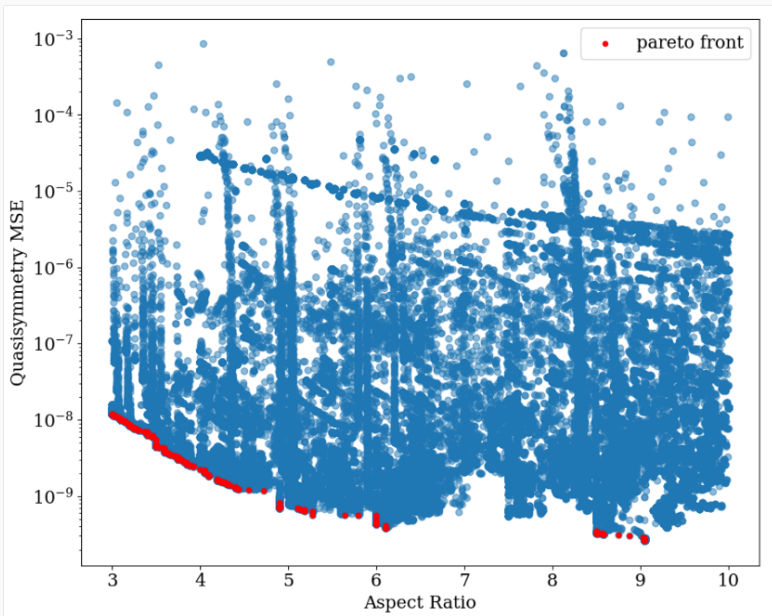
Example: Quasi-symmetry



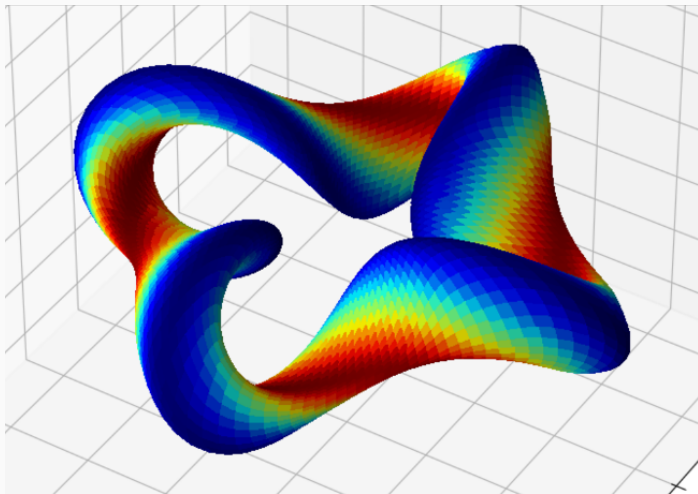
Landreman-Paul QA and QH configurations,
optimized with target aspect ratio 6 and 8.

Q: tradeoff between quasisymmetry and aspect ratio?
(Padidar, Landreman, Bindel)

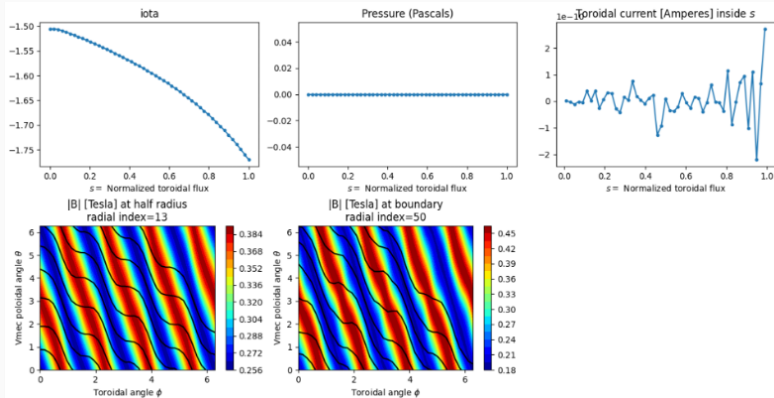
Pareto frontier (QH with 4 field periods)



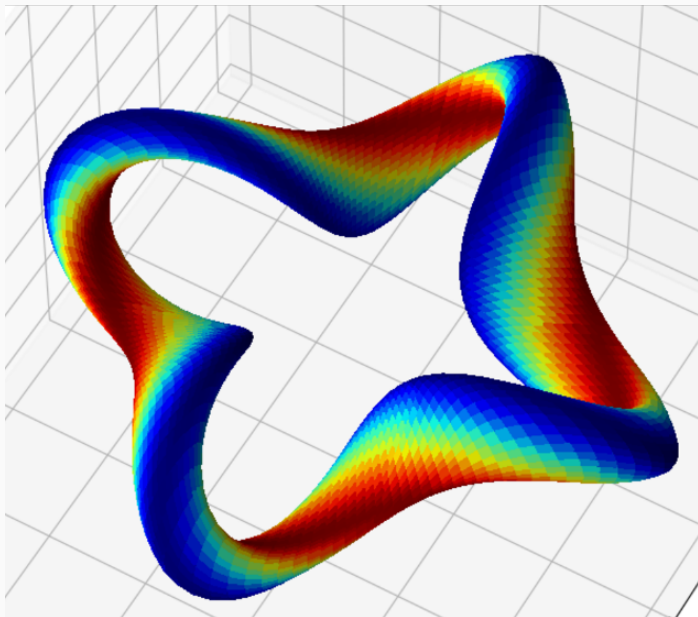
Aspect ratio 3.3



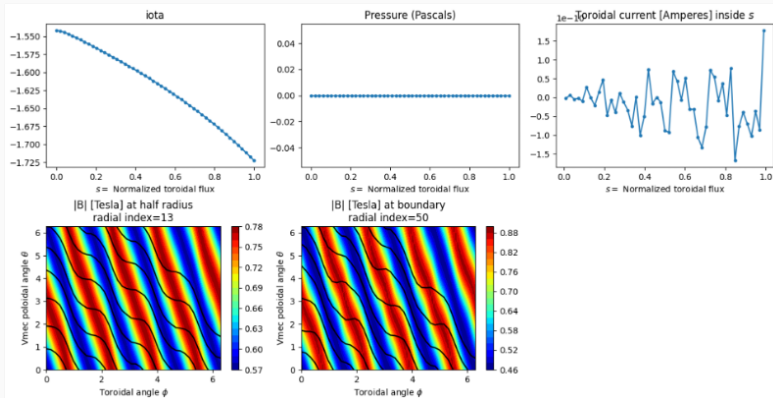
Aspect ratio 3.3



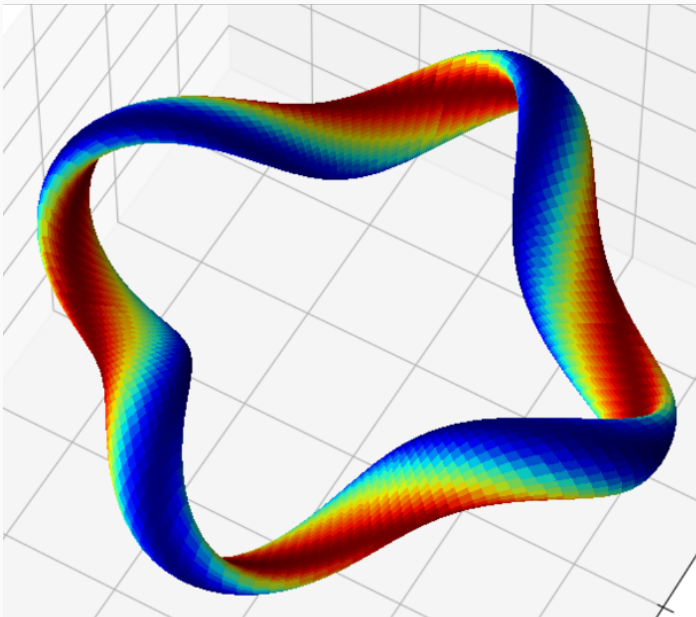
Aspect ratio 5



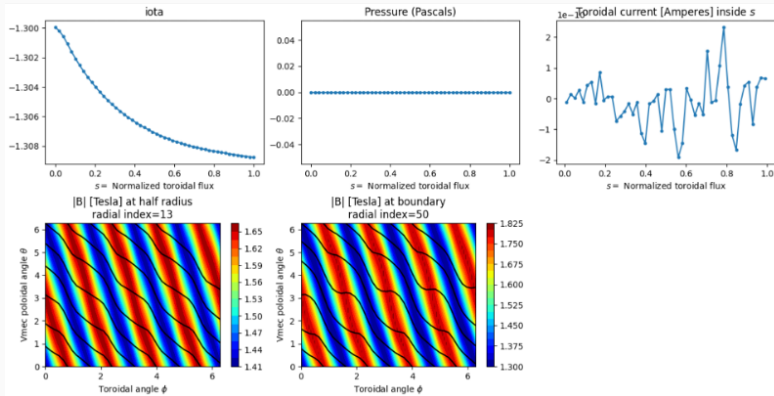
Aspect ratio 5



Aspect ratio 8.67



Aspect ratio 8.67



Continuation

Algorithm in this case: continuation in A

- Start at one Pareto point $(A(x), Q(x))$
- Write stationarity conditions via

$$\nabla Q(x) + \lambda \nabla A(x) = 0$$

$$\lambda(A(x) - A^*) = 0$$

$$A(x) \leq A^*$$

- Differentiate vs A^* to get tangent direction

$$\begin{bmatrix} \nabla^2 Q(x) + \lambda \nabla^2 A(x) & \nabla A(x) \\ \nabla A(x)^T & 0 \end{bmatrix} \begin{bmatrix} x' \\ \lambda' \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

- Predictor moves a little in tangent direction
- Correct prediction via local solver (e.g. Newton)
- Can re-use Hessians, etc for more efficiency

Which parameterization?

What if Pareto frontier goes vertical?

- Can switch to using Q as continuation parameter
- Or use a **pseudo-arclength** parameter
- Generalizations to more than two functions are available (e.g. normal boundary intersection)

And more!

Feel free to ask about

- The latest with GPs with derivatives (Xinran)
- Optimization with stability constraints (Max)
- Fast computations of flux surfaces (Max)
- Alpha particle transport (Misha)