# Data sampling for surrogate modeling and optimization

Tyler Chang
(and others)

Argonne National Laboratory

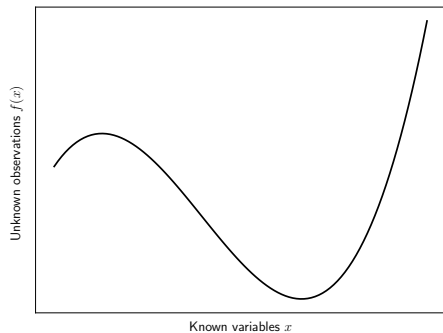ICIAM 2023, Tokyo, Japan
Aug 23, 2023

# Outlines

Inference problems, the curse of dimensionality, and measure collapse

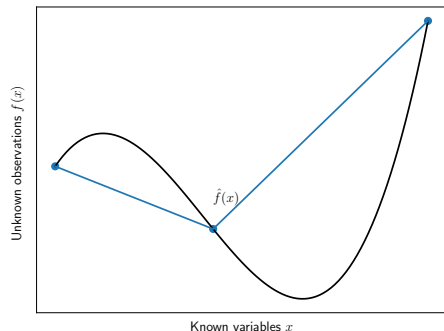Modeling for high-dimensional optimization

Applications

# The fundamental machine learning problem
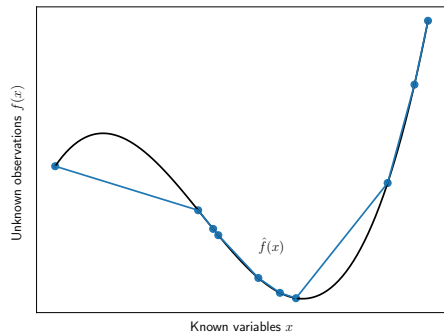
# The fundamental machine learning problem



Known variables $x$ / Unknown observations $f(x)$

▶ Want to predict unknown $f(x)$ for observation $x$

# The fundamental machine learning problem



- ▶ Want to predict unknown $f(x)$ for observation $x$
- ▶ **ML**: *Learn* approximation $\hat{f} \sim f$ based on *training data* $\mathcal{X}$
- ▶ **NA**: fit an interpolant (piecewise-linear) to $f$ on $\mathcal{X}$

# The fundamental machine learning problem



Known variables $x$

(y-axis label: Unknown observations $f(x)$, with $\hat{f}(x)$ marked on the curve)
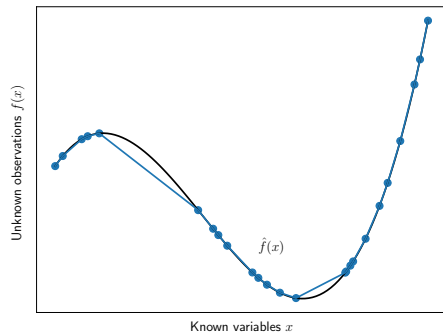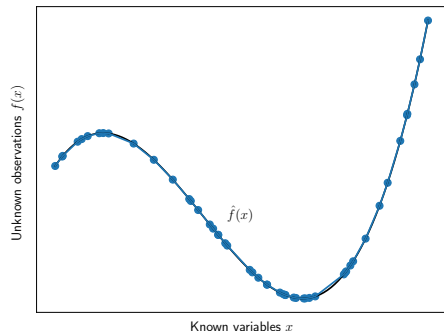
- ▶ Want to predict unknown $f(x)$ for observation $x$
- ▶ **ML**: *Learn* approximation $\hat{f} \sim f$ based on *training data* $\mathcal{X}$
- ▶ **NA**: fit an interpolant (piecewise-linear) to $f$ on $\mathcal{X}$
- ▶ Both cases: more data $\Rightarrow$ better $\hat{f}$

# The fundamental machine learning problem



Known variables $x$ / Unknown observations $f(x)$ / $\hat{f}(x)$

- ▶ Want to predict unknown $f(x)$ for observation $x$
- ▶ **ML**: *Learn* approximation $\hat{f} \sim f$ based on *training data* $\mathcal{X}$
- ▶ **NA**: fit an interpolant (piecewise-linear) to $f$ on $\mathcal{X}$
- ▶ Both cases: more data $\Rightarrow$ better $\hat{f}$
- ▶ Real data not perfectly balanced $\Rightarrow$ $\hat{f} \to f$ non-uniformly

# The fundamental machine learning problem



- ▶ Want to predict unknown $f(x)$ for observation $x$
- ▶ **ML**: *Learn* approximation $\hat{f} \sim f$ based on *training data* $\mathcal{X}$
- ▶ **NA**: fit an interpolant (piecewise-linear) to $f$ on $\mathcal{X}$
- ▶ Both cases: more data $\Rightarrow$ better $\hat{f}$
- ▶ Real data not perfectly balanced $\Rightarrow$ $\hat{f} \to f$ non-uniformly
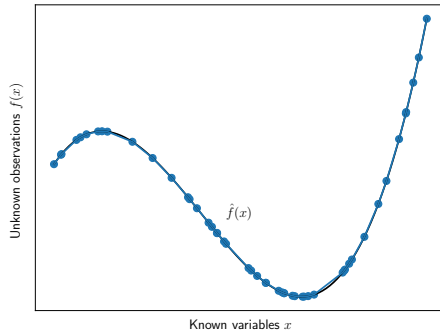- ▶ If we have enough data, it doesn't matter

# Some basic numerical analysis results

When $\hat{f}$ is a piecewise linear spline:

For $h$ "small enough" – let $q$ be the querry point

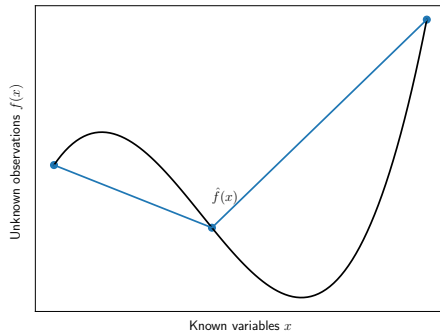$$|f(q) - \hat{f}(q)| \sim \mathcal{O}(h^2)$$



- $h$ is a "mesh fineness" parameter $\sim$ distance between points in $\mathcal{X}$
- For irregular $\mathcal{X}$, $h$ could be the distance from $q$ to the nearest neighbor in $\mathcal{X}$
- Constants proportional to the Lip constant of $\nabla f$

# Some basic numerical analysis results

When $\hat{f}$ is a piecewise linear spline:

For $h$ "small enough" – let $q$ be the querry point

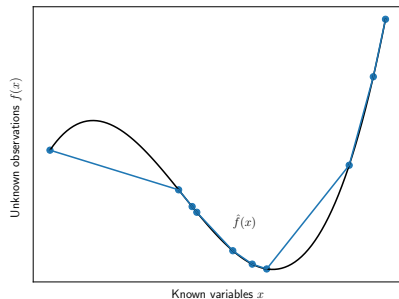$$|f(q) - \hat{f}(q)| \sim \mathcal{O}(h^2)$$



- $h$ is a "mesh fineness" parameter $\sim$ distance between points in $\mathcal{X}$
- For irregular $\mathcal{X}$, $h$ could be the distance from $q$ to the nearest neighbor in $\mathcal{X}$
- Constants proportional to the Lip constant of $\nabla f$
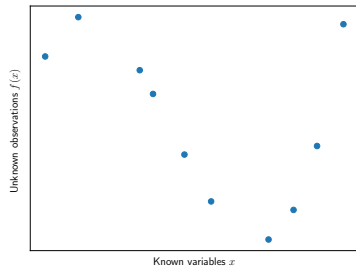
# Some basic deep learning

- ▶ Train a fully-connected multi-layer perceptron (MLP) using $\mathcal{X}$
- ▶ The most popular activation function is ReLU (piecewise linear)
- ▶ In modern ML, train as close to zero error as possible (interpolate)
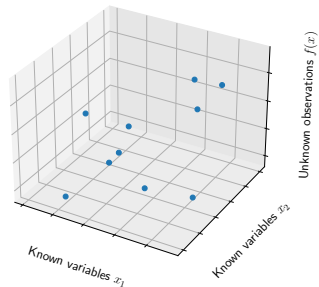
# Some basic deep learning

- ▶ Train a fully-connected multi-layer perceptron (MLP) using $\mathcal{X}$
- ▶ The most popular activation function is ReLU (piecewise linear)
- ▶ In modern ML, train as close to zero error as possible (interpolate)
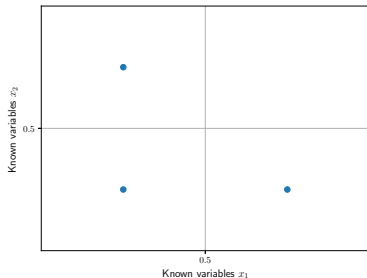
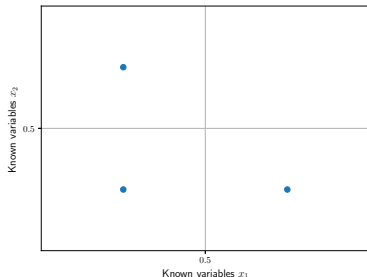# The curse of dimensionality



10 training points in 1D



10 training points in 2D

# The curse of ~~dimensionality~~ no data



Need data in all quadrants?

Need data in all quadrants?

- Inference in 2D : $2^2 = 4$
- Inference in 10D : $2^{10} \approx 1000$
- Inference in 100D : $2^{100} \approx 10^{30}$ (orders of magnitude bigger than exascale)
- Many ML problems : inference in $1000+$ dimensions

# Measure collapse

Can we still make good predictions where we **do** have data?

## Measure collapse

Can we still make good predictions where we **do** have data?

**No, because we have no data anywhere**

We measure where we *might* have enough data to make a prediction using the "convex hull" of the training data $CH(\mathcal{X})$

## Measure collapse

Can we still make good predictions where we **do** have data?

**No, because we have no data anywhere**

We measure where we *might* have enough data to make a prediction using the "convex hull" of the training data $CH(\mathcal{X})$

If $\mathcal{X}$ are sampled from *any* distribution, $\mu(CH(\mathcal{X})) \to 0$ *exponentially* as $d$ grows

This is called a *concentration of measure*

Gorban and Tyukin, Stochastic separation theorems. *Neural Networks 94*, pp. 255-259 (2017).

## Example

Suppose that we uniformly sample $x = (x_1, x_2, \ldots, x_d)$ from $[0,1]^d$
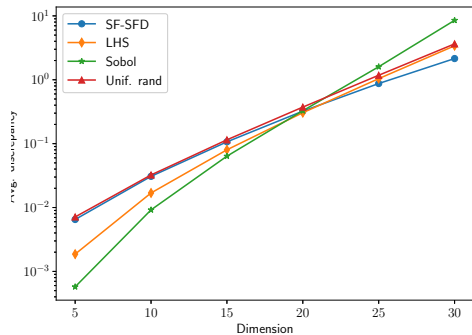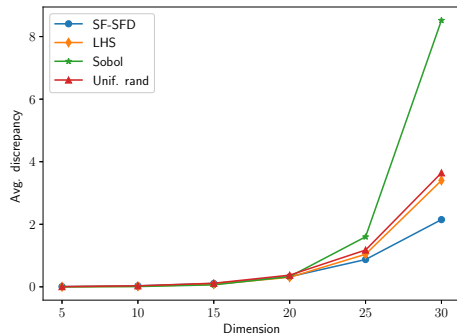
$$\|x - \frac{1}{2}\|_2^2 = \sum_{i=1}^{d} (x_i - \frac{1}{2})^2.$$

$$\mathbb{E}\left[\left(x_i - \frac{1}{2}\right)^2\right] = \int_0^1 \left(u - \frac{1}{2}\right)^2 du = \frac{1}{12}$$

with finite variance $v$

By CLT for all $x \in \mathcal{X}$: $\mathbb{E}[\|x - \frac{1}{2}\|_2^2] = \frac{d}{12}$ with variance $\frac{v}{d} \to 0$ as $d \to \infty$.
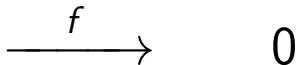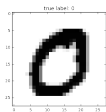
# Collapse of some common distributions



Garg, Chang, and Raghavan, Stochastic optimization of Fourier coefficiencts to generate space-filling designs. *To appear in Winter Sim 2023.*

# Representation learning solution

**"There's more to machine learning than function approximation"**
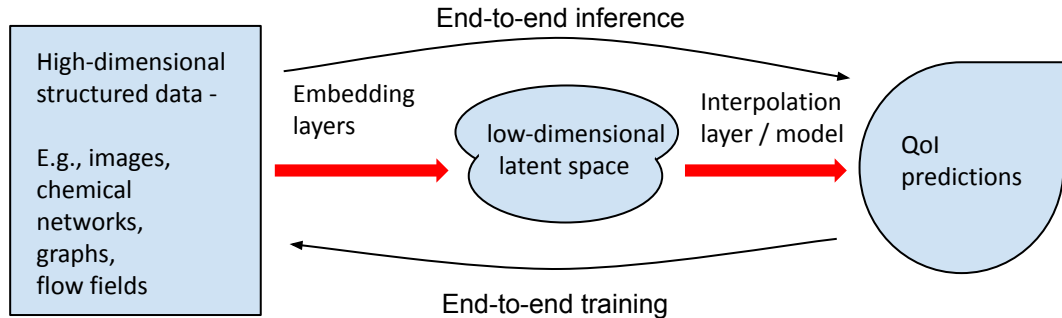
# Representation learning solution

**"There's more to machine learning than function approximation"**

▶ $f$ is often highly *structured* – MLPs with nothing else are from the 60s



$$\xrightarrow{\;f\;}\quad 0$$

$28 \times 28$ pixels $\neq 784$ dimensions...

# Modern deep learning pipeline

End-to-end inference

High-dimensional structured data -

E.g., images, chemical networks, graphs, flow fields

Embedding layers

low-dimensional latent space

Interpolation layer / model

QoI predictions

End-to-end training

# Multiobjective Black-Box Optimization



design space
$\mathcal{X}$

soln set ↗

$\longrightarrow$

Objective space ($\mathcal{Y}$)

Pareto front

# Global optimization

In global optimization literature...

▶ Balance exploration vs. exploitation

▶ Drive *global model error* to zero

▶ Need exponentially many samples to guarantee global convergence

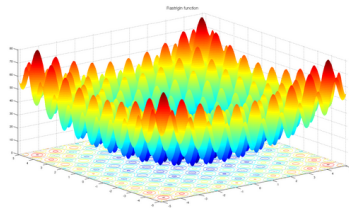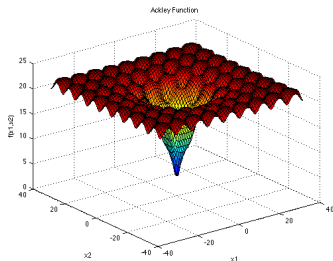Guarantees convergence for problems with thousands of local minima

# Local optimization

▶ Only exploit – maybe multi-start or large initial search
▶ Fit a model that is *locally accurate*
  ▶ Sample requirement grows only *linearly with* dimension
▶ Modification is as simple as putting a *trust-region* around interesting points

# Local optimization

▶ Only exploit – maybe multi-start or large initial search
▶ Fit a model that is *locally accurate*
  ▶ Sample requirement grows only *linearly with* dimension
▶ Modification is as simple as putting a *trust-region* around interesting points

# Local optimization

- ▶ Only exploit – maybe multi-start or large initial search
- ▶ Fit a model that is *locally accurate*
  - ▶ Sample requirement grows only *linearly with* dimension
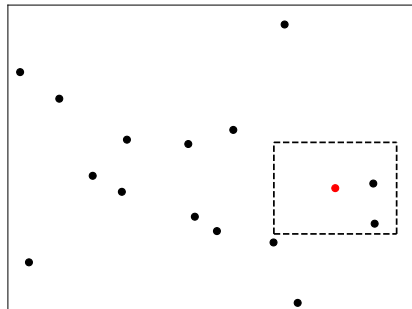- ▶ Modification is as simple as putting a *trust-region* around interesting points

# Local optimization

▶ Only exploit – maybe multi-start or large initial search
▶ Fit a model that is *locally accurate*
   ▶ Sample requirement grows only *linearly with* dimension
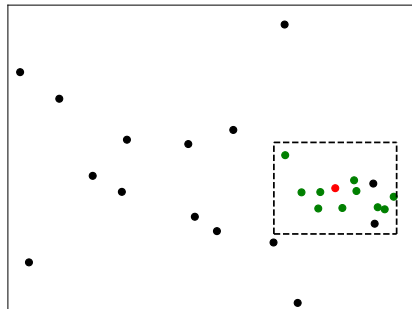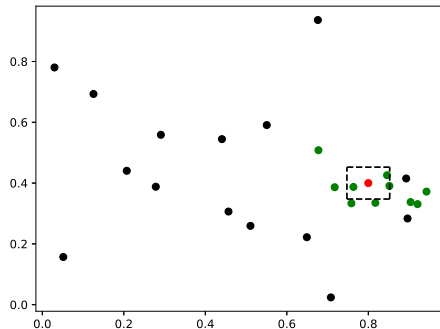▶ Modification is as simple as putting a *trust-region* around interesting points

# ParMOO



Written in `Python`

Version 0.3.0 is now available on available on `pip`,
`conda-forge`, and `GitHub`

`https://github.com/parmoo/parmoo`

`https://parmoo.readthedocs.io`

Chang and Wild. ParMOO: A Python library for parallel multiobjective simulation optimization. *JOSS 8(82):4468 (2023)*.

# Chemical Design on a Limited Budget

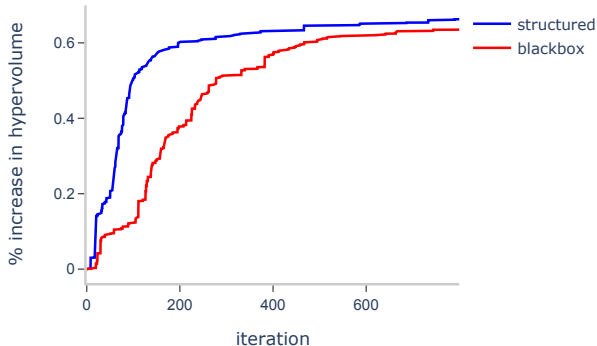- ▶ 6-dimensional latent space embedding of a mixed-variable problem
- ▶ 3-objectives electrolyte manufacturing
  - ▶ high yield, minimal byproduct, low reaction times
- ▶ Running real-world experiments with very limited budget



Chang et al. A framework for fully autonomous design of materials via multiobjective optimization and active learning: challenges and next steps *In ICLR 2023, Workshop on ML4Materials*.
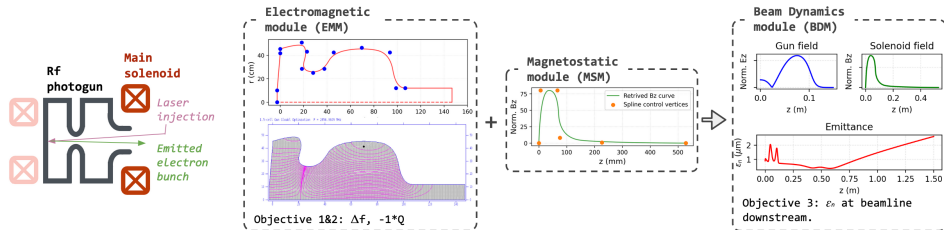
# Fayans Model Calibration (Inverse Problem)

- ▶ 13-variable, 3-objective problem
- ▶ Higher dimensional, requires trust-region methods



Chang and Wild. Designing a framework for solving multiobjective simulation optimization problems. *Under review, preprint* https://arxiv.org/abs/2304.06881.

# Particle Accelerator Beam Design

- ▶ 22-variable, 2-objective problem
- ▶ 3 physics constraints, nearly impossible to satisfy
- ▶ Matched well-known reference gun geometry with just **1300** true simulation evaluations



Chen, Chang, et al. An Integrated Multi-Physics Optimization Framework for Particle Accelerator Design. *Under review.*

## Some Conclusions

- ▶ Doing anything global (modeling, optimization, etc.) in high-dimensions is very hard (maybe impossible)
- ▶ Easier to identify low-dimensional structures and model these *locally*
  - ▶ In my experience, giving up global accuracy is the only thing that scales to big problems
- ▶ Some problems (optimization) don't necessarilly require global accuracy
  - ▶ Don't demand it if you don't need it!
- ▶ Optimization rarely truly requires global accuracy

# Some Conclusions

▶ Doing anything global (modeling, optimization, etc.) in high-dimensions is very hard (maybe impossible)
▶ Easier to identify low-dimensional structures and model these *locally*
  ▶ In my experience, giving up global accuracy is the only thing that scales to big problems
▶ Some problems (optimization) don't necessarilly require global accuracy
  ▶ Don't demand it if you don't need it!
▶ Optimization rarely truly requires global accuracy
▶ *But there are other problems that do require global accuracy...*

# References

*Garg, Chang, and Raghavan. Stochastic optimization of Fourier coefficiencts to generate space-filling designs. To appear in Winter Sim 2023.*

*Chang and Wild. ParMOO: A Python library for parallel multiobjective simulation optimization. JOSS 8(82):4468 (2023).*

*Chang and Wild. Designing a framework for solving multiobjective simulation optimization problems. Under Review, ArXiv preprint 2304.06881 (2023).*

*Chang et al. A framework for fully autonomous design of materials via multiobjective optimization and active learning: challenges and next steps. In ICLR 2023, Workshop on ML4Materials.*

*Chen, Chang, et al. An Integrated Multi-Physics Optimization Framework for Particle Accelerator Design. Under review.*

# Resources

GitHub: github.com/parmoo/parmoo

Pip: pip install parmoo
Conda: conda install --channel=conda-forge parmoo

Test problems: github.com/parmoo/parmoo-solver-farm

tchang@anl.gov