

50ème Journées de Statistique

1 Juin 2018

**Adaptivity and Structure in
Optimization and Estimation**

John Lafferty

Department of Statistics and Data Science
Yale University

Theme

This talk is about the interplay between global and local structure in statistical learning. Two different settings:

- Stochastic convex optimization
- Shape constrained estimation

Outline

- I. Adaptivity and geometry in optimization
 - ▶ How hard is it to optimize *this* function?

- II. Graph structured signal denoising
 - ▶ What is behavior of shape constraints on graphs?

- III. Prediction rule reshaping
 - ▶ How can shape constraints be used with ML algorithms?

I: Geometry and adaptivity in optimization

Heard around the Chicago Statistics lunch table:

“Computer scientists are pessimists”

“I don't care about minimax”

The pessimism is about worst-case thinking, which is overly conservative.

What are alternatives?

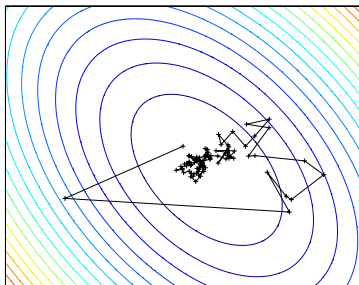
Per-Instance Complexity

$$\begin{aligned} & \text{minimize } f(x) \\ & \text{subject to } x \in \mathcal{C} \end{aligned}$$

f convex, \mathcal{C} closed and convex.

Algorithms get noisy gradients of f at T query points.

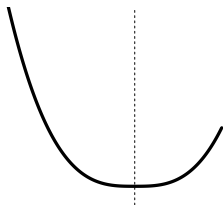
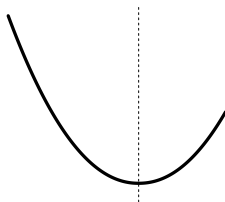
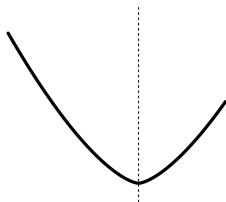
Example: logistic regression or SVMs for large-scale datasets.



Question: *How hard is it to optimize this function f ?*

Per-Instance Complexity

- How hard is it to optimize this function?





Sabyasachi Chatterjee
UIUC



John Duchi
Stanford

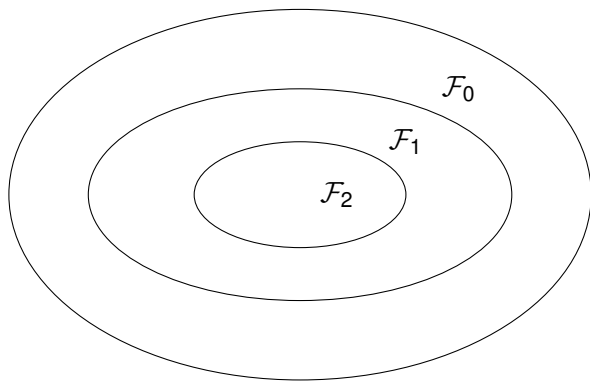


Yuancheng Zhu
UPenn

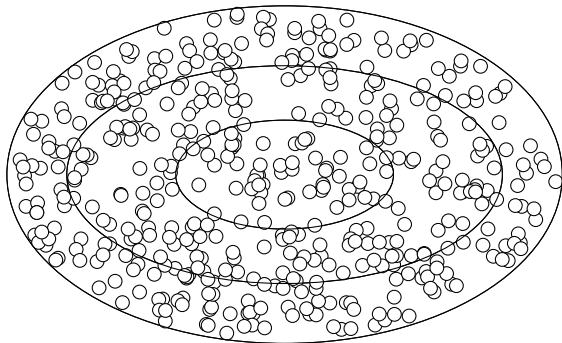
Perspective

- Stochastic gradient descent central to modern machine learning
- Our work shows that SGD is adaptive on a very fine scale
- Making sense of this requires a break with traditional formulations
 - ▶ A complexity class with a single member doesn't make sense using standard notions
- Ideas bridging optimization and statistics

Traditional adaptivity



Superadaptivity



Minimax complexity of convex optimization

- \mathcal{F} *class of convex functions* on a convex set $\mathcal{C} \subset \mathbb{R}^d$.
- \mathcal{O} a *stochastic first-order oracle*: query $(f, x) \in \mathcal{F} \times \mathcal{C}$, returns $Z \in \mathbb{R}^d$, with mean $f'(x) \in \partial f(x)$.
- \mathcal{A}_T class of *all optimization methods* that make T queries to \mathcal{O} .

Minimax complexity

$$R_T(\mathcal{F}) = \inf_{A \in \mathcal{A}_T} \sup_{f \in \mathcal{F}} \mathbb{E}[\text{err}(f, A)] = \inf_{A \in \mathcal{A}_T} \sup_{f \in \mathcal{F}} \mathbb{E} \left[f(x_{T+1}) - \inf_{x \in \mathcal{C}} f(x) \right]$$

Minimax complexity for convex optimization

Known that

$R_T(\mathcal{F}_{\text{sc}}) \asymp 1/T$ strongly convex functions

$R_T(\mathcal{F}_L) \asymp 1/\sqrt{T}$ Lipschitz functions

Agarwal et al. (2010) extend analysis to d -dimensional case, also sparse setting

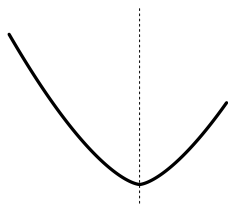
Raginsky and Rakhlin (2011) information theoretic proof technique; parallels minimax lower bounds in statistics

Shortcomings of the framework?

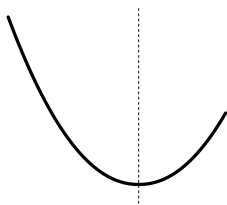
- Ignores cost of computing the gradient – $O(1)$
- Does not allow for (decreasing) bias
- Does not account for computations on past gradients, e.g., quasi-Newton algorithms
- *Too pessimistic and “global”*

Per-instance complexity

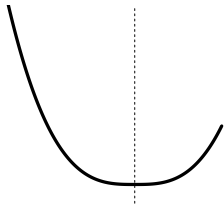
What about a particular instance?



$$\ll 1/T?$$



$$\asymp 1/T?$$



$$\ll 1/\sqrt{T}?$$

Local minimax complexity

Minimax complexity

$$R_T(\mathcal{F}) = \inf_{A \in \mathcal{A}_T} \sup_{f \in \mathcal{F}} \mathbb{E}[\text{err}(f, A)]$$

Local minimax complexity

$$R_T(f; \mathcal{F}) = \sup_{g \in \mathcal{F}} \inf_{A \in \mathcal{A}_T} \max_{h \in \{f, g\}} \mathbb{E}[\text{err}(h, A)]$$

Why is this interesting?

Local minimax complexity

$$R_T(f; \mathcal{F}) = \sup_{g \in \mathcal{F}} \inf_{A \in \mathcal{A}_T} \max_{h \in \{f, g\}} \mathbb{E}[\text{err}(h, A)]$$

To show this is a reasonable complexity measure, need to:

1. Relate it to geometry of the function
2. Show it agrees with rates for known classes
3. Demonstrate that beating it is impossible
4. Give an algorithm that achieves it

Definitions

Set of minimizers $\mathcal{X}_f^* = \arg \min_{x \in \mathcal{C}} f(x)$

Error function $\text{err}(x, f) = \inf_{y \in \mathcal{X}_f^*} \|x - y\|$

Minima separation $d(f, g) = \inf_{x \in \mathcal{X}_f^*, y \in \mathcal{X}_g^*} \|x - y\|$ for $f, g \in \mathcal{F}$

Subgradient gap $\kappa(f, g) = \sup_{x \in \mathcal{C}} \|f'(x) - g'(x)\|$

Modulus of continuity

Two dissimilarity measures between f and g :

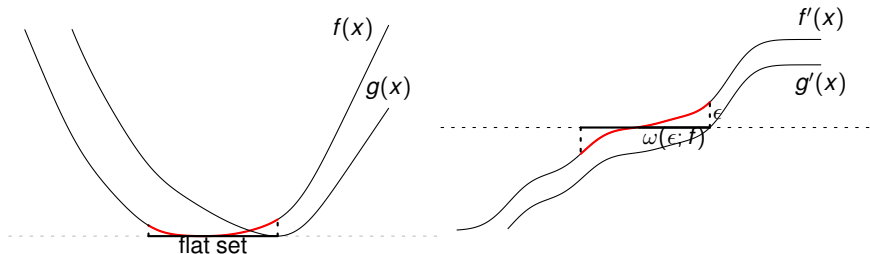
$d(f, g)$ *distance between minimizers*

$\kappa(f, g)$ *largest separation between subgradients*

Modulus of continuity of d with respect to κ at function f :

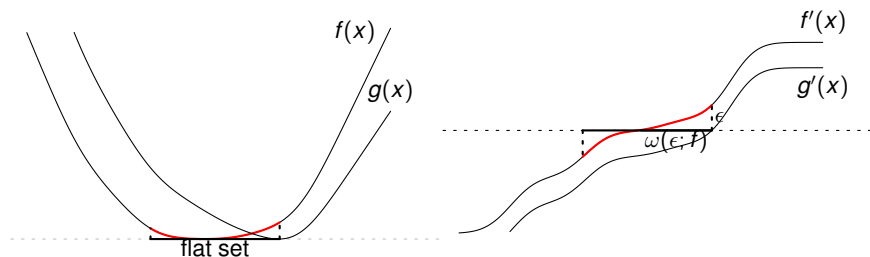
$$\omega_f(\epsilon) = \sup_{g \in \mathcal{F}} \{d(f, g) : \kappa(f, g) \leq \epsilon\}$$

Modulus of continuity



$$\omega_f(\epsilon) = \sup \left\{ \inf_{x \in \mathcal{X}_f^*} |x - y| : y \in \mathcal{C}, |f'(y)| < \epsilon \right\}$$

Modulus of continuity

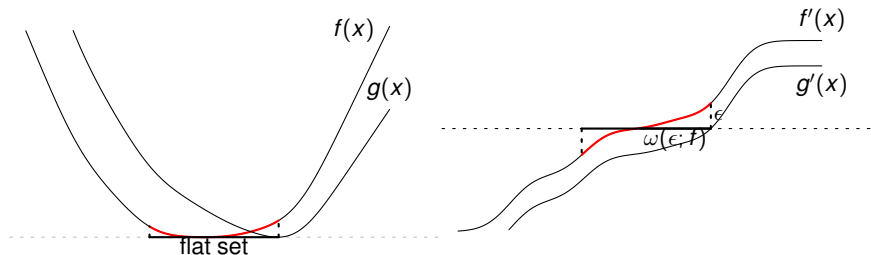


$$\omega_f(\epsilon) = \sup \left\{ \inf_{x \in \mathcal{X}_f^*} |x - y| : y \in \mathcal{C}, |f'(y)| < \epsilon \right\}$$

$$f(x) = \frac{1}{\alpha} |x|^\alpha$$

$$\omega_f(\epsilon) = \epsilon^{1/(\alpha-1)}$$

Modulus of continuity



$$\omega_f(\epsilon) = \sup \left\{ \inf_{x \in \mathcal{X}_f^*} |x - y| : y \in \mathcal{C}, |f'(y)| < \epsilon \right\}$$

$$f(x) = \frac{1}{\alpha} |x|^\alpha$$

$$\omega_f(\epsilon) = \epsilon^{1/(\alpha-1)}$$

$$f(x) = \begin{cases} \frac{1}{\alpha} |x|^\alpha & x \leq 0 \\ \frac{1}{\beta} |x|^\beta & x > 0 \end{cases}$$

$$\omega_f(\epsilon) = \epsilon^{1/(\alpha \vee \beta - 1)}$$

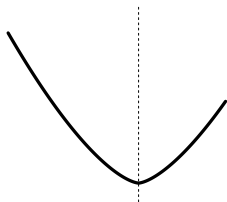
Modulus characterizes local minimax

Theorem. For all sufficiently large T ,

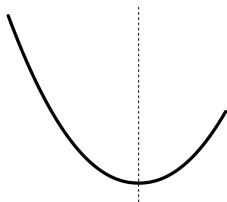
$$C_1 \omega_f \left(\frac{\sigma}{\sqrt{T}} \right) \leq R_T(f; \mathcal{F}) \leq C_2 \omega_f \left(\frac{\sigma}{\sqrt{T}} \right).$$

- Proof works for several distances $d(f, g)$
- For $f(x) = c|x - x^*|^\alpha$, x -error decays as $O(T^{-1/2(\alpha-1)})$,
 f -error decays as $O(T^{-\alpha/2(\alpha-1)})$
- Agrees with known rates for uniformly convex functions

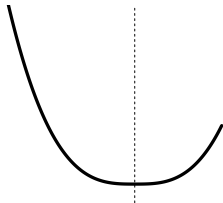
Per-instance complexity



$$\asymp 1/T^{3/2}$$



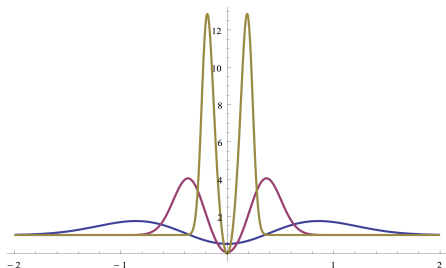
$$\asymp 1/T$$



$$\asymp 1/T^{2/3}$$

Superefficiency

- In statistics, Fisher information is the benchmark for efficiency of (parametric) estimators
- The MLE satisfies $\sqrt{n}(\hat{\theta}_n - \theta) \rightsquigarrow N(0, I(\theta)^{-1})$
- LeCam-Hájek: If an estimator satisfies $\sqrt{n}(\tilde{\theta}_n - \theta) \rightsquigarrow N(0, v(\theta))$ with $v(\theta) \ll 1/I(\theta)$, then it must perform poorly at a nearby point



Superefficiency in convex optimization

Suppose an algorithm $A \in \mathcal{A}_T$ outperforms the modulus:

$$\mathbb{E}_f \text{err}(\hat{x}_A, f) \leq \delta_T \omega_f \left(\frac{\sigma}{\sqrt{T}} \right),$$

with $\delta_T \rightarrow 0$, $e^T \delta_T \rightarrow \infty$. Then exists functions with $\kappa(f, g_T) \rightarrow 0$ and

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{g_T} \text{err}(\hat{x}_A, g_T)}{\omega_{g_T} \left(\sigma \sqrt{T^{-1} \log(1/\delta_T)} \right)} > 0.$$

Thus ω_f can be viewed as analogue of Fisher information for stochastic convex optimization

A superadaptive algorithm: Binary search

Given a budget of T queries:

- Query midpoint $T / \log T$ times
- Average the noisy gradients
- If result is positive, go left; otherwise right
- Step is correct if derivative bigger than C/\sqrt{T} .

Binary search algorithm

Input: T, r

Initialize: (a_0, b_0) , $E = \lfloor r \log T \rfloor$, $T_0 = \lfloor T/E \rfloor$,

for $e = 1, \dots, E$ **do:**

Query $x_e = (a_e + b_e)/2$ for T_0 times to get $Z_t^{(e)}$ for $t = 1, \dots, T_0$

Calculate the average $\bar{Z}_{T_0}^{(e)} = \frac{1}{T_0} \sum_{t=1}^{T_0} Z_t^{(e)}$

If $\bar{Z}_{T_0}^{(e)} > 0$, set $(a_{e+1}, b_{e+1}) = (a_e, x_e)$

If $\bar{Z}_{T_0}^{(e)} \leq 0$, set $(a_{e+1}, b_{e+1}) = (x_e, b_e)$

end

Output: x_E

Binary search achieves the benchmark

Theorem. With probability at least $1 - \delta$ and for large enough T ,

$$\inf_{x \in \mathcal{X}_f^*} |x_T - x| \leq \tilde{C} \omega_f \left(\frac{\sigma}{\sqrt{T}} \right)$$

where the term \tilde{C} hides a dependence on $\log T$ and $\log(1/\delta)$.

Blockwise SGD is superadaptive in $d > 1$

- Can't (easily) do binary search in higher dimension
- Nesterov's blockwise SGD is extremely simple:
 - 1 Divide computational budget into equal chunks
 - 2 Run SGD in each block with fixed step size
 - 3 Halve the step size in each block
- Superadaptive (up to log factors)

Other settings

Recent results of Yuancheng Zhu

	zeroth-order oracle	first-order oracle
fixed/random design	regression	1st-order regression
active design	0th-order optimization	optimization

$$R_{*,\dagger}(n, f) = \sup_{g \in \mathcal{F}} \inf_{A \in \mathcal{A}_{n,\dagger}} \max_{h \in \{f, g\}} \mathbb{E}_h \text{err}_*(A, h)$$

See also...

J. Duchi and F. Ruan (2017), “Local asymptotics for stochastic optimization: Optimality, constraint identification, and dual averaging”

- Local minimax for optimization analogous to Hájek-Le Cam
- Based on Nesterov’s dual averaging

W. Su and Y. Zhu (2018), “Statistical inference for online learning and stochastic approximation via hierarchical incremental gradient descent”

- Inference (confidence intervals) for SGD
- Hierarchical design, computationally “free”

Part I: Summary

- Framework for assessing complexity of minimizing individual convex functions.
- Close connections to new and old statistical theory
- Challenge: More natural definition?

II: Graph structured signals

- Estimation and testing of signals that respect the structure of a network or graph in some way.
- We formulate a form of isotonic regression on graphs, and study the risk properties of the least squares estimator.

Trend filtering on graphs, Wang et al., (2014)

Normal means on graphs, Arias-Castro et al., (2008, 2014), Sharpnack (2013)

Lipschitz learning on graphs, Kyng et al., (2015)

II: Graph structured signals

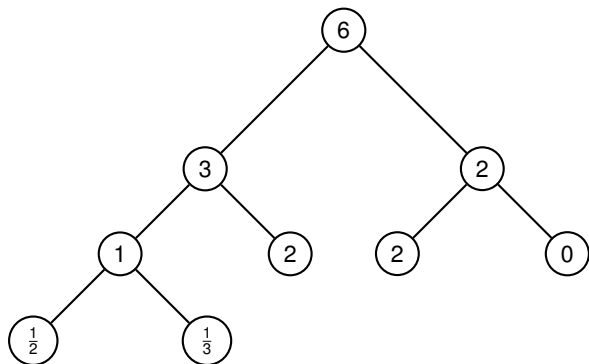


Sabyasachi Chatterjee

UIUC

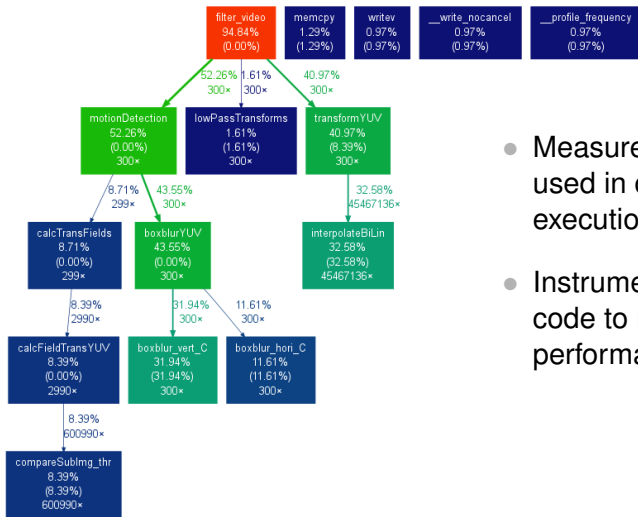
Flows on trees

Imagine a fluid flowing into a node and dividing among the children – possibly with some leakage.



We observe a noisy measurement $X_i = \mu_i + \epsilon_i$ at each node i

Example: Statistical code profiling (gprof)



- Measures time/storage used in different parts of execution tree
- Instruments compiled code to monitor performance

Denoising flows

- Graph structured form of isotonic regression.
 - ▶ What is the behavior of the least-squares estimator?
 - ▶ How does it depend on the structure of the tree?
 - ▶ What is the fundamental limit of flow estimation?

The isotonic case

For isotonic regression, $\mu_1 \geq \dots \geq \mu_n$, risk of LSE scales as

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}(\mu_i - \hat{\mu}_i)^2 \leq C \left(\frac{\sigma^2(\mu_1 - \mu_n)}{n} \right)^{2/3} = O(n^{-2/3})$$

Matches the minimax rate.

LSE for tree flows: Logarithmic depth

Theorem. Let T_n be a sequence of trees with n nodes and depth h_n . For any flow $\mu \in \mathcal{F}(T_n)$, the LSE has risk

$$\frac{1}{n} \mathbb{E} \|\hat{\mu} - \mu\|^2 \leq C \left(\frac{\sigma^2 h_n (1 + \log n)^3}{n} + \frac{\sigma \mu_1 \sqrt{h_n} (1 + \log n)^{3/2}}{n} \right)$$

where μ_1 is the flow at the root.

LSE for tree flows: Logarithmic depth

Theorem. Let T_n be a sequence of trees with n nodes and depth h_n . For any flow $\mu \in \mathcal{F}(T_n)$, the LSE has risk

$$\begin{aligned} \frac{1}{n} \mathbb{E} \|\hat{\mu} - \mu\|^2 &\leq C \left(\frac{\sigma^2 h_n (1 + \log n)^3}{n} + \frac{\sigma \mu_1 \sqrt{h_n} (1 + \log n)^{3/2}}{n} \right) \\ &= \tilde{O} \left(\frac{h_n}{n} \right) \end{aligned}$$

where μ_1 is the flow at the root.

LSE for tree flows: Logarithmic depth

Theorem. Let T_n be a sequence of trees with n nodes and depth h_n . For any flow $\mu \in \mathcal{F}(T_n)$, the LSE has risk

$$\begin{aligned} \frac{1}{n} \mathbb{E} \|\hat{\mu} - \mu\|^2 &\leq C \left(\frac{\sigma^2 h_n (1 + \log n)^3}{n} + \frac{\sigma \mu_1 \sqrt{h_n} (1 + \log n)^{3/2}}{n} \right) \\ &= \tilde{O} \left(\frac{h_n}{n} + \frac{\mu_1 \sqrt{h_n}}{n} \right) \end{aligned}$$

where μ_1 is the flow at the root.

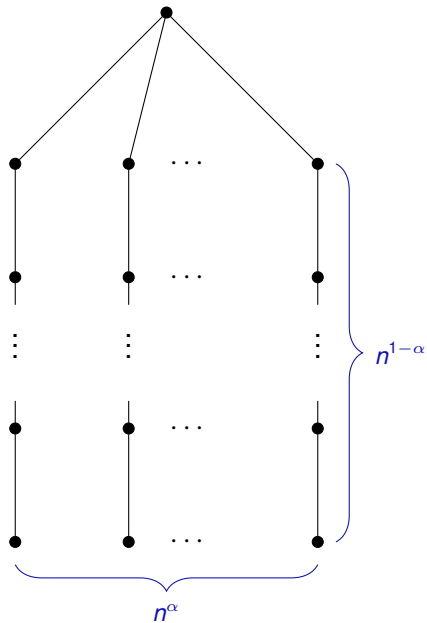
Hardest trees?

- Flow estimation is easier for stars than for paths

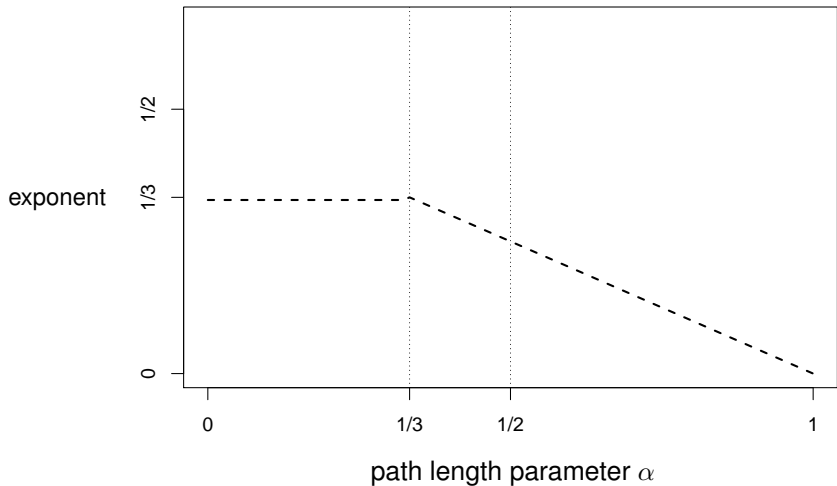
$$\tilde{O}(1/n) \text{ vs. } O(n^{-2/3})$$

- Is the path the “hardest” flow estimation problem?

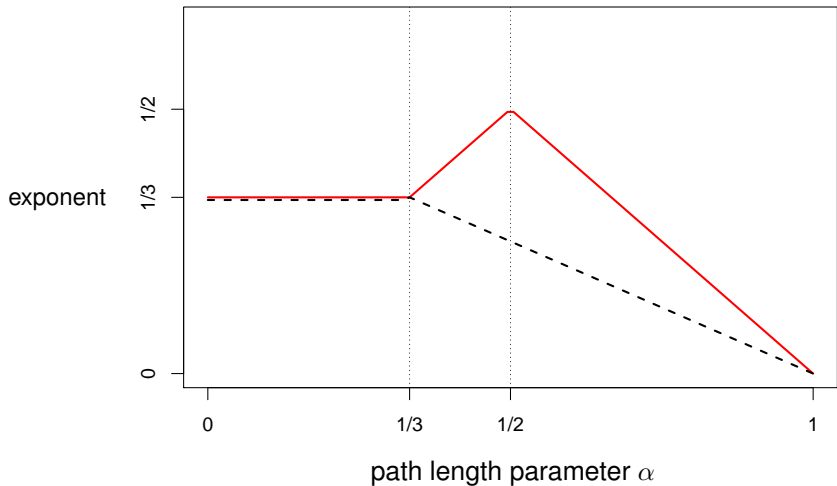
Many long paths $\mathcal{T}_{n,\alpha}$



Lower bounds for $\mathcal{T}_{n,\alpha}$



Lower bounds for $\mathcal{T}_{n,\alpha}$ and **LSE**



Proof techniques

<i>Result</i>	<i>Technique</i>
Upper bound for LSE; shallow trees	Gaussian supremum
Minimax lower bound; shallow trees	Fano's lemma
Lower bound for LSE; shallow trees	Gaussian widths
Isotonic upper bound for LSE; deep trees	Statistical dimension
Simplex upper bound for LSE; deep trees	Chaining
Minimax lower bound; monotone seqs	Assouad's lemma
Minimax lower bound; deep trees, $\alpha \leq \frac{1}{3}$	Minimax for isotonic
Minimax lower bound; deep trees, $\alpha \geq \frac{1}{3}$	Fano's lemma
Tightness of lower bound, $\alpha \geq \frac{1}{3}$	LSE on net
Tightness of LSE upper bound, $\alpha \geq \frac{1}{3}$	Gaussian widths

Proof techniques: GP suprema

For a fixed flow μ we define the Gaussian supremum function

$$f_{\mu}(t) := \mathbb{E} \left(\sup_{\nu \in \mathcal{F}: \|\nu - \mu\| \leq t} \langle Z, \nu - \mu \rangle \right) - \frac{t^2}{2}.$$

If $t^* > 0$ satisfies $f_{\mu}(t^*) \leq 0$ then

$$R(\hat{\mu}, \mu) \leq \frac{C}{n} \max(t^{*2}, \sigma^2).$$

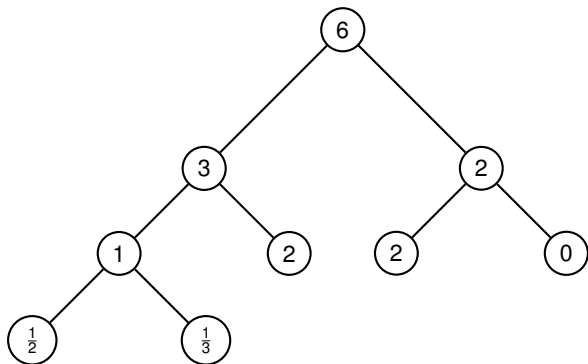
We bound $f_{\mu}(t)$ by Dudley's entropy integrals (chaining)

$$\mathbb{E} \left(\sup_{\nu \in \mathcal{F}: \|\nu - \mu\| \leq t} \langle Z, \nu - \mu \rangle \right) \leq C \int_0^{2t} \sqrt{\log N(\epsilon, B_t(\mu))} d\epsilon$$

which requires good upper bounds on log covering numbers.

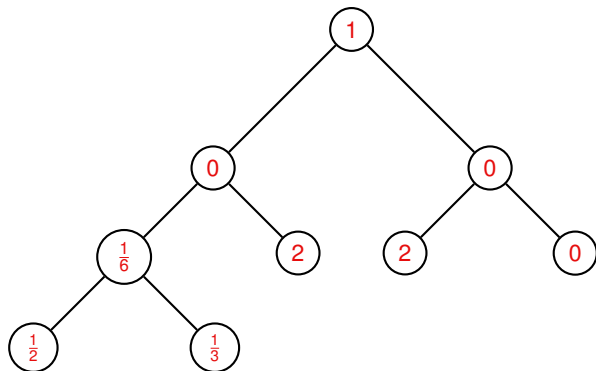
Covering number sketch: Flows and leaks

A flow is determined by its leaks. Let $\mu_1 \leq V$.



Covering number sketch: Flows and leaks

A flow is determined by its leaks. Let $\mu_1 \leq V$.



Random construction

Given a flow μ , define a random flow F by its leaks:

$$\ell(F) = \begin{cases} Ve_j & \text{with probability } \frac{\ell_j(\mu)}{V} \\ 0 & \text{with probability } 1 - \frac{1}{V} \sum_i \ell_i(\mu). \end{cases}$$

Take convex combination of m random flows:

$$\bar{\mu}^m = \frac{1}{m}(F_1 + \dots + F_m)$$

Recursion

Since the node flows are negatively correlated, can show

$$\text{Var}(\bar{\mu}_v^m) \leq \sum_{u \in \text{Subtree}(v)} \text{Var}(\bar{\ell}_u^m)$$

$$\begin{aligned} \mathbb{E} \|\bar{\mu}^m - \mu\|^2 &= \sum_i \mathbb{E} (\bar{\mu}_i^m - \mu_i)^2 = \sum_i \text{Var}(\bar{\mu}_i^m) \\ &\leq \sum_i d_i \text{Var}(\bar{\ell}_i^m) \\ &\leq h \sum_i \text{Var}(\bar{\ell}_i^m) \\ &\leq \frac{V^2 h}{m} \end{aligned}$$

Covering number

Random construction thus gives a $V^2 h/m$ covering. By simple combinatorics,

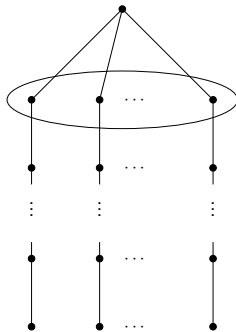
$$\log N(\epsilon, \mathcal{F}_V) \leq \frac{V^2 h}{\epsilon^2} \left(1 + \log \left(1 + \frac{n \epsilon^2}{V^2 h} \right) \right)$$

Intuition for the gap

General minimax lower bound for isotonic regression:

$$\inf_{\tilde{\mu}} \sup_{\mu: \mu_1 \leq V} \mathbb{E} \|\tilde{\mu} - \mu\|^2 \geq C \min \left\{ \sigma^2, V^2, \left(\frac{\sigma^2 V}{n} \right)^{2/3} \right\}$$

- When V is small, LSE not minimax—trivial estimator $\hat{\mu} = 0$ achieves lower bound
- In “narrow” tree regime $\alpha \approx 1/3$, some of the root flows will be small
- A lasso-style estimator to zero entire paths?



II: Summary

- Flows: Tree structured generalization of isotonic regression
- Surprise: LSE is not minimax rate optimal
- Shape constraints in graph/network settings largely unexplored

S. Chatterjee and JL, “Denoising flows on trees,” *IEEE Trans. Info. Theory*, 2018

III: Prediction rule reshaping

- Imposing shape constraints on “black box” prediction rules
- Reshaping random forests



Matt Bonakdarpour
UChicago/Yale



Rina Foygel Barber
UChicago

Motivation

- Shape constraints are natural in applications
 - ▶ House price assumed to be decreasing function of crime rate (all other predictors held constant)
- Not easily incorporated into popular machine algorithms
- We study different approaches to *reshaping* prediction rules

Prediction rule reshaping

Suppose that $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ is a prediction rule estimated from data (e.g., regression or classification)

Let \mathcal{S} be a class of functions that satisfy a set of shape constraints. For example

$$\frac{\partial f}{\partial x_j} \geq 0 \quad (\text{monotonic}), j \in \mathcal{R}_\uparrow$$
$$\frac{\partial^2 f}{\partial x_j^2} \leq 0 \quad (\text{concave}), j \in \mathcal{R}_\cap$$

Reshaping is the infinite dimensional optimization

$$\min_{f \in \mathcal{S}} \|f - \hat{f}\|$$

A finite optimization problem

Let $\mathcal{S} = (\mathcal{S}_v)_{v \in \mathcal{R}}$ be candidate shape constraints, $\mathcal{R} \subset [d]$

Let $\mathcal{D}_n = \{x_1, \dots, x_n\}$ be set of test points

For $v \in \mathcal{R}$, define the $n \times n$ prediction matrix $\widehat{F}^v = [\widehat{F}_{i',j}^v]$

$$\widehat{F}_{i',j}^v = \widehat{f}(x_{i,1}, x_{i,2}, \dots, x_{i,v-1}, x_{i',v}, x_{i,v+1}, \dots, x_{i,d})$$

$$\widehat{F}_i^v(\cdot) \equiv \widehat{f}(x_{i,1}, x_{i,2}, \dots, x_{i,v-1}, \cdot, x_{i,v+1}, \dots, x_{i,d}).$$

where $x_{i',v}$ ranges over all n values of v -th predictor.

A finite optimization problem

Reshaped predictions \check{f}_S obtained by projecting matrix of predicted values onto shape constraints:

$$\check{F}_S = \arg \min_{F=(F^v)_{v \in \mathcal{R}}} \sum_{v \in \mathcal{R}} \|F^v - \hat{F}^v\|_F^2$$

such that $F_i^v \in \mathcal{S}_v$, for each $v \in \mathcal{R}$
 $\text{diag}(F^v) = \text{diag}(F^w)$, for all $v, w \in \mathcal{R}$

Reshaped predictions:

$$\check{f}_S = \text{diag}(\check{F}_S^v)$$

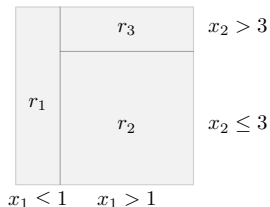
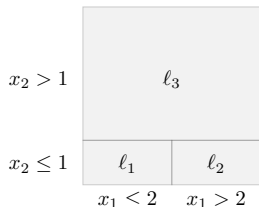
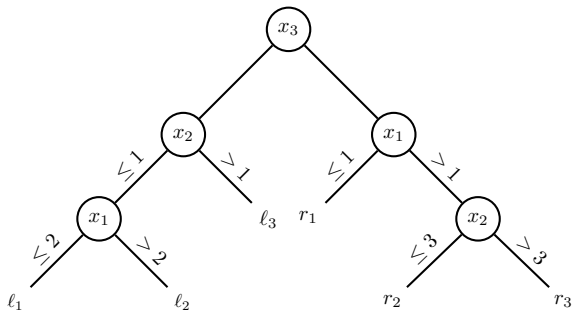
Intersecting isotonic regressions

- Suppose constraints are monotonicity constraints
- Reshaping leads to “intersecting isotonic regressions” due to consistency constraints $\text{diag}(F^V) = \text{diag}(F^W)$
- A generalization of PAVA solves this efficiently

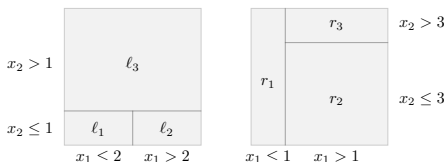
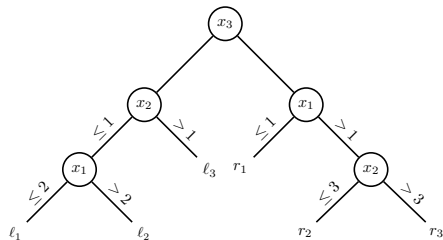
Reshaping random forests

- 1 Grow the tree in the usual way
- 2 Reshape the leaf values to enforce monotonicity

Illustration



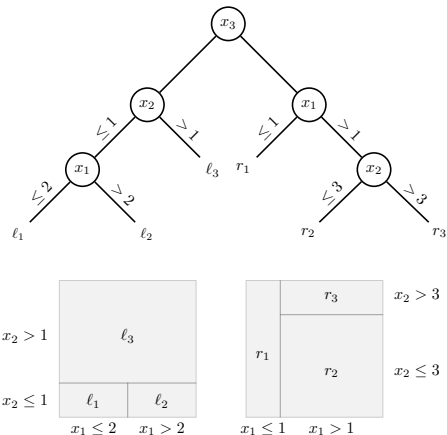
Illustration



Exact estimator, 6 constraints:

$$\mu_{l_2} \leq \mu_{r_2}, \mu_{l_1} \leq \mu_{r_1}, \mu_{l_1} \leq \mu_{r_2}, \mu_{l_3} \leq \mu_{r_1}, \mu_{l_3} \leq \mu_{r_2}, \mu_{l_3} \leq \mu_{r_3}$$

Illustration



Over-constrained estimator, all 9 pairwise constraints

$$\mu_{\ell_1} \leq \mu_{r_1}, \mu_{\ell_1} \leq \mu_{r_2}, \mu_{\ell_1} \leq \mu_{r_3}, \dots$$

Reference

“Prediction rule reshaping,” Matt Bonakdarpour, Sabyasachi Chatterjee, Rina Foygel Barber, and JL, arXiv:1805.06439

(to be presented at ICML this summer)

Summary

Interplay between global and local structure:

- Stochastic convex optimization
- Signal denoising on trees
- Reshaping prediction rules

Merci beaucoup d'avoir m'écouté!