Parallel Computing on Big Data

- In the parallel computing environment, a common practice is to distribute a massive dataset to multiple processors, and then aggregate local results obtained from separate machines into global counterparts;
- The above Divide-and-Conquer (D&C) strategy often requires a growing number of machines to deal with an increasingly large dataset;
- This computational consideration leads to the emergence of the so-called "Splitotics Theory," a type of "Big Data Theory."

A Basic Question from Statisticians

- How many machines do we really need in parallel computing from a statistical theory perspective?
- In this poster, we address this basic, yet fundamentally important, question by carefully analyzing statistical versus computational trade-off in the above D&C framework;
- In particular, an intriguing phase transition phenomenon is discovered for the number of deployed machines that ends up being a simple proxy for computing cost, for both statistical estimation and testing.

Divide and Conquer Strategy

A Flowchart of D&C

Nonparametric (univariate) regression model:

$$Y = f_0(Z) + \epsilon.$$
Subset 1 (n) $\stackrel{\text{Machine 1}}{\longrightarrow} \qquad \widehat{f_1}$
Big Data (N) $\stackrel{\text{Divide}}{\longrightarrow} \qquad \begin{array}{c} \text{Subset 2 (n)} \stackrel{\text{Machine 2}}{\longrightarrow} \qquad \widehat{f_2} \\ \dots & \dots & \dots \\ \text{Subset 2 (n)} \stackrel{\text{Machine 2}}{\longrightarrow} \qquad \widehat{f_3} \\ \dots & \dots & \dots \\ \text{Subset s (n)} \stackrel{\text{Machine s}}{\longrightarrow} \qquad \widehat{f_3} \\ \begin{array}{c} \text{Con } \\ \text{Quer} \\ \text{Oracle Est. } \\ \hline{f_N} \\ \end{array} \qquad \begin{array}{c} \text{Agg. Est. } \\ \hline{f_N} \\ \end{array}$

$$\overline{f}_N = (1/s) \sum_{j=1}^s \widehat{f}_n^{(j)}.$$

A Plot for Computing Time

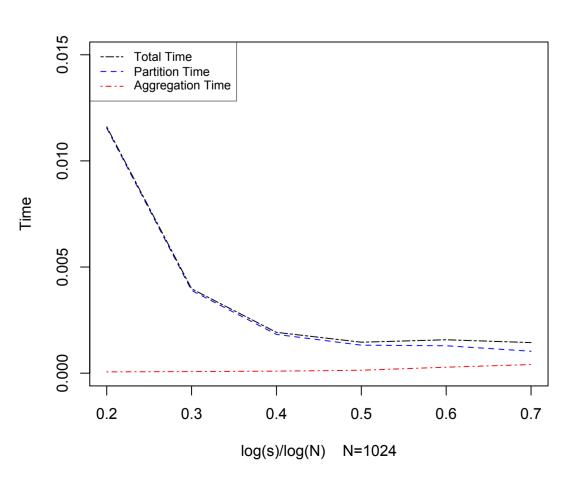


Figure: Computing time for \overline{f}_N when N = 1024 as s diverges

How Many Processors Do We Really Need in Parallel Computing?

Guang Cheng and Zuofeng Shang

Purdue University and Binghamton University

Statistical-and-Computational Tradeoff

A Plot for Mean Squared Error

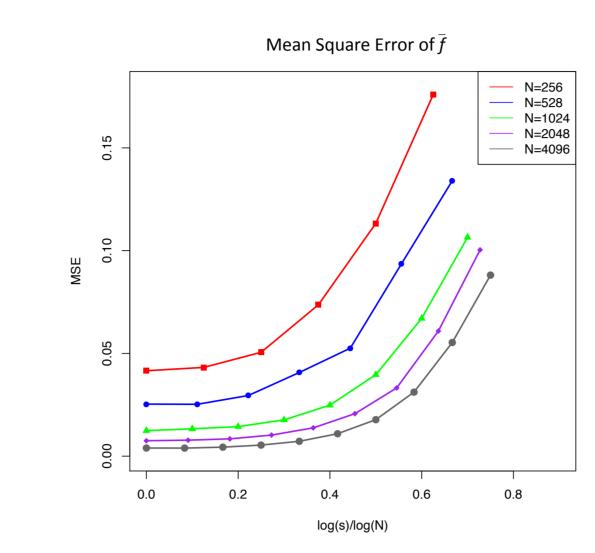


Figure: Mean-square errors of \overline{f}_N under different N as s diverges

Our Major Goal

- We would like to know how fast s is allowed to diverge (w.r.t. N), say $s = N^a$, such that the aggregated estimate \overline{f}_N is minimax optimal or nonparametric testing based on f_N is minimax optimal;
- statistical optimality is achievable and above which statistical optimality is impossible. The sharpness is important in that it captures the *intrinsic* computational limit of the D&C algorithm.

Computational Limit I: Statistical Estimation

Observe samples from the following model

$$\mathbf{y}_{l} = \mathbf{f}(l/N) + \epsilon_{l},$$

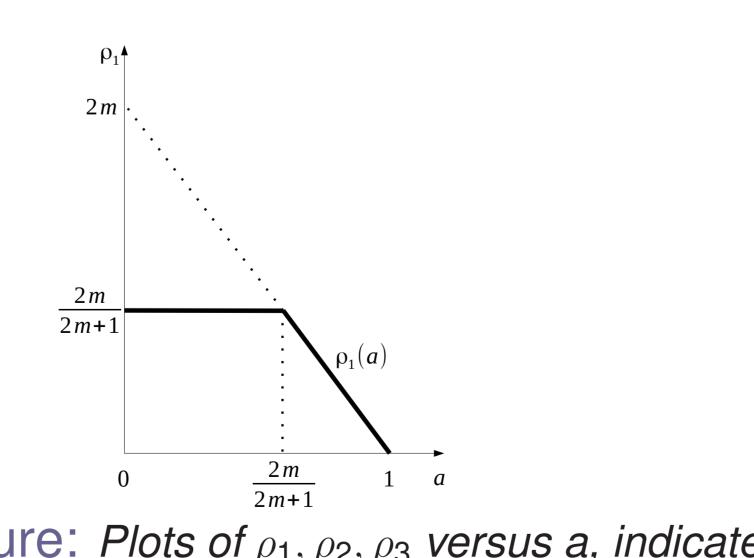
where ϵ_I 's are *iid* zero-mean r.v.s with unit variances;

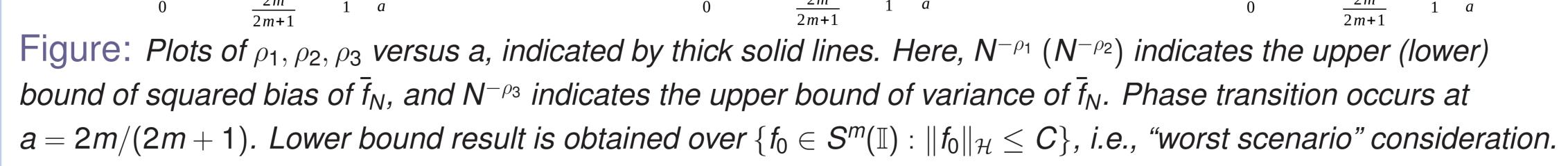
- machines as "evenly" as possible over the entire interval [0, 1];
- At each machine, we obtain a smoothing spline sub-estimate as

$$\widehat{f}_{j} = rg\min_{f\in\mathcal{S}^{m}(\mathbb{I})}rac{1}{n}\sum_{i=1}^{n}(Y_{i,j}-f(t_{i,j}))^{2}+\lambda\|f\|_{\mathcal{H}}^{2},$$

where $\langle f, g \rangle$ is a roughness penalty and $\lambda > 0$ is a smoothing parameter; • $MSE_{f_0}(\bar{f}_N) = E\{\|\bar{f}_N - E\{\bar{f}_N\}\|_2^2\} + \|E\{\bar{f}_N\} - f_0\|_2^2$.

A Graphical Illustration



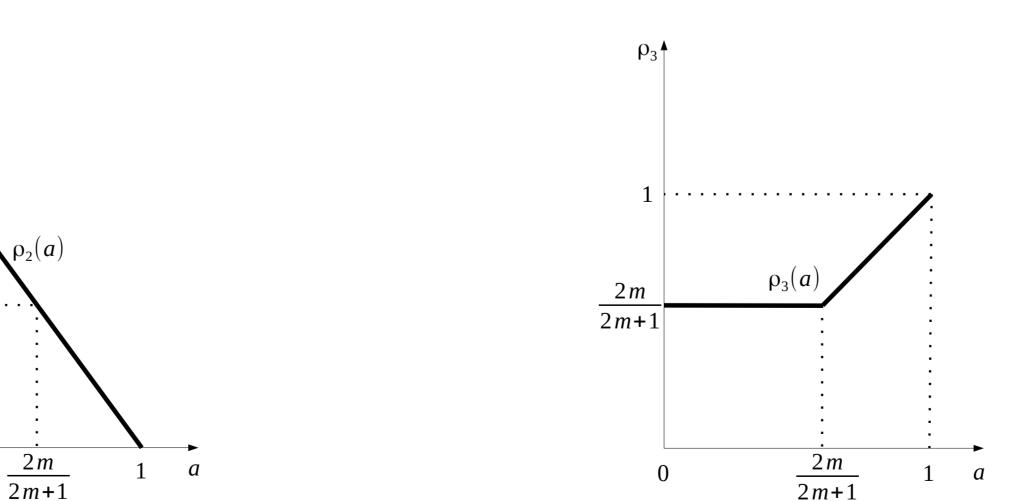


A sharp upper bound $s^* \simeq N^{2m/(2m+1)}$ is established for optimal estimation.

We will show that there indeed exists a sharp upper bound for s, below which

 $I = 0, 1, \ldots, N - 1,$

► The N samples $\{y_l, I/N\}_{l=0}^{N-1}$ are distributed to s machines with each machine being assigned n samples. We want the N covariates $t_l = I/N$ to appear in s

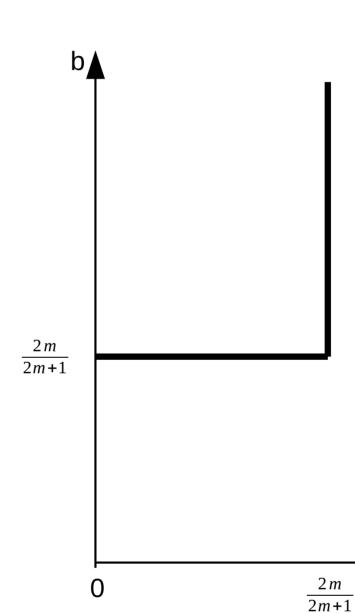


Computational Limit II: Statistical Testing

- Define a Wald-type test statistic:
- Our testing rule is thus

Future Works

Conjecture: is D&C a new form of tuning?



A "Theoretical" Suggestion

Distribute to

- Distribute to

machines for performing an optimal test.

Practical Formulae to be Developed...

method does not work here;

How to select λ that is scaled to N in practice?

Test the following simple hypothesis:

 $H_0: f = 0$ v.s. $H_1: f \in S^m(\mathbb{I}) \setminus \{0\};$

$$T_{N,\lambda} = \|ar{f}_N\|_2^2;$$

 $\phi_{\mathbf{N},\lambda} = I(|T_{\mathbf{N},\lambda} - \mu_{\mathbf{N},\lambda}| \geq Z_{1-\alpha/2}\sigma_{\mathbf{N},\lambda}),$

where $\mu_{N,\lambda}$ and $\sigma_{N,\lambda}^2$ are mean and variance of $T_{N,\lambda}$, respectively;

Testing consistency (Type I error) essentially requires no condition on s as long as $N \to \infty$. In other words, s can be either fixed or diverge at any rate; However, the (non-asymptotic) power of our proposed test depends on s in a very subtle manner. Specifically, the separation rate achieves its minimal value $N^{-2m/(4m+1)}$ under the following sharp upper bound

 $S^{**} \simeq N^{(4m-1)/(4m+1)}$

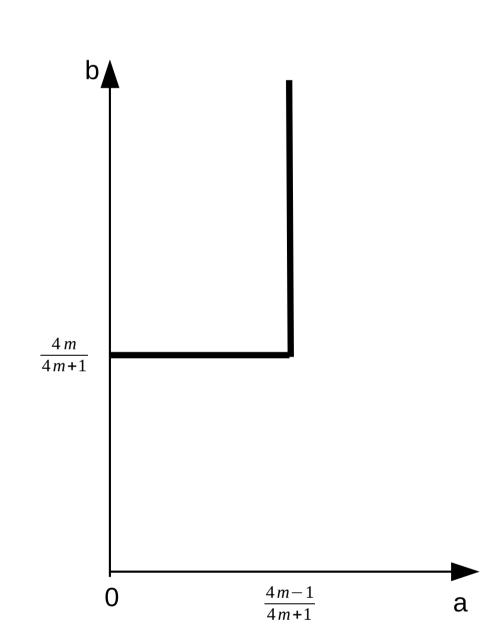


Figure: Two lines indicate the choices of $s \simeq N^a$ and $\lambda \simeq N^{-b}$, leading to minimax optimal estimation (left) and minimax optimal testing (right). Whereas (a, b)'s outside these two lines lead to suboptimal rates.

When applying D&C to massive data, we may allocate machines as:

 $m{s} symp N^{2m/(2m+1)}$

machines for obtaining an optimal estimate;

 $\boldsymbol{S} symp N^{4m/(4m+1)}$

• We show that λ should be chosen in the order of N even when each subsample has size *n*. Hence, the standard generalized cross validation

So far, we can only give a theoretical upper bound for s.

How to pick the number of machines in practice using data-depdent method?