

Personalized Treatment: Sounds heavenly, but where on Earth did they find *my* guinea pigs?

Xiao-Li Meng, Harvard University

- Meng (2014) A trio of inference problems that could win you a Nobel Prize in statistics (if you help fund it). *In Past, Present and Future of Statistics*. COPSS 50th Anniversary Vol. (Ed: Lin et. al)
- Liu and Meng (2014). A fruitful resolution to Simpson's paradox via multi-resolution inference. *The American Statistician*, 68: 17-29.
 - Liu and Meng (2016). There is individualized treatment. Why not
- individualized inference?. Annual Review of Statistics and Its Application, 3, 79-111.

A Multi-resolution Theory for Approximating Infinite-*p*-Zero-*n*: Transitional Inference, Individualized Predictions, and a World without Bias-variance Trade-off (*JASA* (2021) special issue on "Precision medicine and individualized policy discovery")

Xinran Li and Xiao-Li Meng

- Meng (2014) A trio of inference problems that could win you a Nobel Prize in statistics (if you help fund it). In Past, Present and Future of Statistics. COPSS 50th Anniversary Volume (Ed: Lin et. al.)
- Liu and Meng (2014). A fruitful resolution to Simpson's paradox via multi-resolution inference. *The American Statistician*, 68: 17-29.
- Liu and Meng (2016). There is individualized treatment. Why not individualized inference?. Annual Review of Statistics and Its Application, 3, 79-111.

Xinran Li and Xiao-Li Meng

Э

イロト イポト イヨト イヨト

A Big Picture about Data Science

"Whereas previous generations were in possession of data about general phenomena, we are now in possession of data about specific phenomena. For example, in genomics we have data about each individual gene, in astronomy we have data about each region of the sky, in medicine we have data about each tumor, and in social science we have data about individual humans. Our era is about 'data' and about 'specific context.' In that sense 'data science' is an appropriate and useful terminology for capturing current trends."

> Ani Adhikari, John DeNero, and Michael Jordan (*Harvard Data Science Review*, Winter, 2021)

Treatment for you based only on data from people like you.

Э

イロト イポト イヨト イヨト

Treatment for you based only on data from people like you. BUT: No one is perfectly like you.



Treatment for you based only on data from people like you.

BUT: No one is perfectly like you.



Relevance for COVID-19: Population Risk vs Individual Risk

Xinran Li and Xiao-Li Meng

э

ヨト イヨト

- Potentially Infinitely many attributes: $p = \infty$
- Each of us is unique: direct learning sample size n = 0.



- Potentially Infinitely many attributes: $p = \infty$
- Each of us is unique: direct learning sample size n = 0.



 Moving from population-level "soft matching" to individual-level "hard matching" – is that possible?

- Potentially Infinitely many attributes: $p = \infty$
- Each of us is unique: direct learning sample size n = 0.



• Moving from population-level "soft matching" to individual-level "hard matching" – is that possible?

Ultimate bias-variance trade-off

Relevance vs Robustness

• "Transition to the similar"

э

프 에 에 프 어

Image: Image:

• "Transition to the similar"

Aelius Galenus (129-210 AD)

- ∢ ∃ →

• "Transition to the similar"

Aelius Galenus (129-210 AD)

• Transitional inference: an empiricism concept (Hankinson 1987, 1995)



• "Transition to the similar"

Aelius Galenus (129-210 AD)

• Transitional inference: an empiricism concept (Hankinson 1987, 1995)



• "In cases in which there is no history, or in which there is none of sufficient similarity, there is not much hope.

• "Transition to the similar"

Aelius Galenus (129-210 AD)

• Transitional inference: an empiricism concept (Hankinson 1987, 1995)



• "In cases in which there is no history, or in which there is none of sufficient similarity, there is not much hope. And the same thing is true in the case of transference of one remedy from one ailment to another similar to it: one has a greater or smaller basis for expectation of success in proportion to the increase or decrease in similarity of the ailment, whether or not history is involved.

• "Transition to the similar"

Aelius Galenus (129-210 AD)

• Transitional inference: an empiricism concept (Hankinson 1987, 1995)



• "In cases in which there is no history, or in which there is none of sufficient similarity, there is not much hope. And the same thing is true in the case of transference of one remedy from one ailment to another similar to it: one has a greater or smaller basis for expectation of success in proportion to the increase or decrease in similarity of the ailment, whether or not history is involved. And the same goes for the transference from one part of the body to another part: expectation of success varies in direct proportion to the similarity."

Inference/Prediction as Approximation

• What Do We Want to Know?: Your response to treatment t, $Y_t(\odot^*)$, for t = A and t = B.

6

▶ ∢ ∃ ▶

Inference/Prediction as Approximation

- What Do We Want to Know?: Your response to treatment t, $Y_t(\odot^*)$, for t = A and t = B.
- What Do We Know?: Either Y_A(☺) or Y_B(☺) (but not both) for some other people, ☺ ≠ ☺*.

Inference/Prediction as Approximation

- What Do We Want to Know?: Your response to treatment t, $Y_t(\odot^*)$, for t = A and t = B.
- What Do We Know?: Either Y_A(☺) or Y_B(☺) (but not both) for some other people, ☺ ≠ ☺*.
- **Strategy**: Must construct population of "relevant" individuals with which to approximate **YOU**.

• The ©* relevant subpopulation is

$$\Omega_{\mathcal{C}}(\odot^*) = \{ \odot : \mathcal{C}(\odot) = \mathcal{C}(\odot^*) \}.$$

• The ©* relevant subpopulation is

$$\Omega_{\mathcal{C}}(\odot^*) = \{ \odot : \mathcal{C}(\odot) = \mathcal{C}(\odot^*) \}.$$

where C is a set of intrinsic (pre-treatment) characteristics.

• $R = \dim(C)$ is the primary resolution.

• The ©* relevant subpopulation is

$$\Omega_{\mathcal{C}}(\odot^*) = \{ \odot : \mathcal{C}(\odot) = \mathcal{C}(\odot^*) \}.$$

- $R = \dim(C)$ is the primary resolution.
- Inspired by multi-resolution wavelets formulation: match on "signals" (low resolutions) and ignore "noises" (high resolutions)

• The ©* relevant subpopulation is

$$\Omega_{\mathcal{C}}(\odot^*) = \{ \odot : \mathcal{C}(\odot) = \mathcal{C}(\odot^*) \}.$$

- $R = \dim(C)$ is the primary resolution.
- Inspired by multi-resolution wavelets formulation: match on "signals" (low resolutions) and ignore "noises" (high resolutions)

• The ©* relevant subpopulation is

$$\Omega_{\mathcal{C}}(\odot^*) = \{ \odot : \mathcal{C}(\odot) = \mathcal{C}(\odot^*) \}.$$

- $R = \dim(C)$ is the primary resolution.
- Inspired by multi-resolution wavelets formulation: match on "signals" (low resolutions) and ignore "noises" (high resolutions) YOU
 A Relevant Individual



A Multi-Resolution View of Big Data

Population Resolution



A Multi-Resolution View of Big Data

Population Resolution

Individual Resolution





How to Capture Resolution?

• Y lives on the same probability space as an information filtration $\{\mathcal{F}_r, r = 0, 1, ..., \}$, e.g., $\mathcal{F}_r = \sigma(X_0, X_1, ..., X_r)$.

9

イヨトイヨト

How to Capture Resolution?

- Y lives on the same probability space as an information filtration $\{\mathcal{F}_r, r = 0, 1, ..., \}$, e.g., $\mathcal{F}_r = \sigma(X_0, X_1, ..., X_r)$.
- r is the index of resolution.

9

ヨト・モート

How to Capture Resolution?

- Y lives on the same probability space as an information filtration $\{\mathcal{F}_r, r = 0, 1, ..., \}$, e.g., $\mathcal{F}_r = \sigma(X_0, X_1, ..., X_r)$.
- *r* is the index of resolution.



9

$$\sigma_r^2 = \mathsf{E}[\sigma_s^2 | \mathcal{F}_r] + \mathsf{E}[(\mu_s - \mu_r)^2 | \mathcal{F}_r]$$

10

∃ ▶ ∢

$$\sigma_r^2 = \mathsf{E}[\sigma_s^2 | \mathcal{F}_r] + \mathsf{E}[(\mu_s - \mu_r)^2 | \mathcal{F}_r]$$

• $(\mu_s - \mu_r)^2$: magnitude of signal at resolution *s*, not modelled at resolution *r*.

$$\sigma_r^2 = \mathsf{E}[\sigma_s^2 | \mathcal{F}_r] + \mathsf{E}[(\mu_s - \mu_r)^2 | \mathcal{F}_r]$$

• $(\mu_s - \mu_r)^2$: magnitude of signal at resolution *s*, not modelled at resolution *r*.



KEY: No such thing called "Variance" except at resolution ∞

$$\sigma_r^2 = \mathsf{E}[\sigma_s^2 | \mathcal{F}_r] + \mathsf{E}[(\mu_s - \mu_r)^2 | \mathcal{F}_r]$$

• $(\mu_s - \mu_r)^2$: magnitude of signal at resolution *s*, not modelled at resolution *r*.



KEY: No such thing called "Variance" except at resolution ∞

(1) If we believe the world is stochastic, then $\sigma_{\infty}^2 > 0$.

$$\sigma_r^2 = \mathsf{E}[\sigma_s^2 | \mathcal{F}_r] + \mathsf{E}[(\mu_s - \mu_r)^2 | \mathcal{F}_r]$$

• $(\mu_s - \mu_r)^2$: magnitude of signal at resolution *s*, not modelled at resolution *r*.



KEY: No such thing called "Variance" except at resolution ∞

(1) If we believe the world is stochastic, then σ²_∞ > 0.
(2) If we believe the world is deterministic/chaotic, then σ²_∞ = 0.

$$\sigma_r^2 = \mathsf{E}[\sigma_s^2 | \mathcal{F}_r] + \mathsf{E}[(\mu_s - \mu_r)^2 | \mathcal{F}_r]$$

• $(\mu_s - \mu_r)^2$: magnitude of signal at resolution *s*, not modelled at resolution *r*.



KEY: No such thing called "Variance" except at resolution ∞

- (1) If we believe the world is stochastic, then $\sigma_{\infty}^2 > 0$.
- (2) If we believe the world is deterministic/chaotic, then $\sigma_{\infty}^2 = 0$.
- (3) There is no general empirical test to distinguish between (1) and (2), with any finite amount of data.
Shifting from $n \rightarrow 0$ to $r \rightarrow \infty$...

• Without direct data, we estimate Y_{me} using $\hat{\mu}_r$ with $r = \tilde{r}$

Errors in Indirect-Data Prediction

$$Y_{me} - \hat{\mu}_{\tilde{r}} = \underbrace{(\mu_{\tilde{r}} - \hat{\mu}_{\tilde{r}})}_{\text{Estimation Error}} + \underbrace{(\mu_{\infty} - \mu_{\tilde{r}})}_{\text{Resolution Error}} + \underbrace{(Y_{me} - \mu_{\infty})}_{\text{Intrinsic Error}}$$

• Without direct data, we estimate Y_{me} using $\hat{\mu}_r$ with $r = \tilde{r}$

Errors in Indirect-Data Prediction

$$Y_{me} - \hat{\mu}_{\tilde{r}} = \underbrace{(\mu_{\tilde{r}} - \hat{\mu}_{\tilde{r}})}_{\text{Estimation Error}} + \underbrace{(\mu_{\infty} - \mu_{\tilde{r}})}_{\text{Resolution Error}} + \underbrace{(Y_{me} - \mu_{\infty})}_{\text{Intrinsic Error}}$$

• Double accent: "hat" for estimation and "tilde" for selection.

• Without direct data, we estimate Y_{me} using $\hat{\mu}_r$ with $r = \tilde{r}$

Errors in Indirect-Data Prediction

$$Y_{me} - \hat{\mu}_{\tilde{r}} = \underbrace{(\mu_{\tilde{r}} - \hat{\mu}_{\tilde{r}})}_{\text{Estimation Error}} + \underbrace{(\mu_{\infty} - \mu_{\tilde{r}})}_{\text{Resolution Error}} + \underbrace{(Y_{me} - \mu_{\infty})}_{\text{Intrinsic Error}}$$

- Double accent: "hat" for estimation and "tilde" for selection.
- How big should \tilde{r} be? A Holy Grail of Statistical Inference and Prediction.

• Without direct data, we estimate Y_{me} using $\hat{\mu}_r$ with $r = \tilde{r}$

Errors in Indirect-Data Prediction

$$Y_{me} - \hat{\mu}_{\tilde{r}} = \underbrace{(\mu_{\tilde{r}} - \hat{\mu}_{\tilde{r}})}_{\text{Estimation Error}} + \underbrace{(\mu_{\infty} - \mu_{\tilde{r}})}_{\text{Resolution Error}} + \underbrace{(Y_{me} - \mu_{\infty})}_{\text{Intrinsic Error}}$$

- Double accent: "hat" for estimation and "tilde" for selection.
- How big should \tilde{r} be? A Holy Grail of Statistical Inference and Prediction.
- Equivalent to bias-variance trade-off, but the resolution framework leads to some surprises.

• Without direct data, we estimate Y_{me} using $\hat{\mu}_r$ with $r = \tilde{r}$

Errors in Indirect-Data Prediction

$$Y_{me} - \hat{\mu}_{\tilde{r}} = \underbrace{(\mu_{\tilde{r}} - \hat{\mu}_{\tilde{r}})}_{\text{Estimation Error}} + \underbrace{(\mu_{\infty} - \mu_{\tilde{r}})}_{\text{Resolution Error}} + \underbrace{(Y_{me} - \mu_{\infty})}_{\text{Intrinsic Error}}$$

- Double accent: "hat" for estimation and "tilde" for selection.
- How big should \tilde{r} be? A Holy Grail of Statistical Inference and Prediction.
- Equivalent to bias-variance trade-off, but the resolution framework leads to some surprises.
- It has the same mathematical setup as **sieve method for non-parametric estimation**.

$$egin{aligned} Y &= \sum_{r=0}^\infty eta_r X_r + \epsilon, \quad \epsilon \sim \mathcal{N}(0, au^2), \ \ \epsilon oxplus ec{X}_\infty, \ \ \mathrm{V}(Y) < \infty. \ X_0 &= 1, \ \ \{X_1, X_2, \ldots\} \ ext{are jointly normal distributed}, \end{aligned}$$

• Assume both target and training populations are generated from

$$egin{aligned} Y &= \sum_{r=0}^\infty eta_r X_r + \epsilon, \quad \epsilon \sim \mathcal{N}(0, au^2), \ \epsilon oxdots oldsymbol{\vec{X}}_\infty, \ \mathrm{V}(Y) < \infty. \ X_0 &= 1, \ \{X_1, X_2, \ldots\} \ ext{are jointly normal distributed}, \end{aligned}$$

• The total prediction error under OLS has three parts:

$$egin{aligned} Y &= \sum_{r=0}^\infty eta_r X_r + \epsilon, \quad \epsilon \sim \mathcal{N}(0, au^2), \ \epsilon oxdots oldsymbol{\vec{X}}_\infty, \ \mathrm{V}(Y) < \infty. \ X_0 &= 1, \ \{X_1, X_2, \ldots\} \ ext{are jointly normal distributed}, \end{aligned}$$

- The total prediction error under OLS has three parts:
 - τ^2 : intrinsic error

$$egin{aligned} Y &= \sum_{r=0}^\infty eta_r X_r + \epsilon, \quad \epsilon \sim \mathcal{N}(0, au^2), \ \epsilon oxdots oldsymbol{\vec{X}}_\infty, \ \mathrm{V}(Y) < \infty. \ X_0 &= 1, \ \{X_1, X_2, \ldots\} \ ext{are jointly normal distributed}, \end{aligned}$$

- The total prediction error under OLS has three parts:
 - τ^2 : intrinsic error
 - $A(r) = \sum_{k=r+1}^{\infty} \Delta_k^2$: approx error, $\Delta_k^2 = V(Y|\vec{X}_{r-1}) V(Y|\vec{X}_r)$

$$egin{aligned} Y &= \sum_{r=0}^\infty eta_r X_r + \epsilon, \quad \epsilon \sim \mathcal{N}(0, au^2), \ \epsilon oxdots ec{X}_\infty, \ \mathrm{V}(Y) < \infty. \ X_0 &= 1, \ \{X_1, X_2, \ldots\} \ ext{are jointly normal distributed}, \end{aligned}$$

- The total prediction error under OLS has three parts:
 - τ^2 : intrinsic error
 - $A(r) = \sum_{k=r+1}^{\infty} \Delta_k^2$: approx error, $\Delta_k^2 = V(Y|\vec{X}_{r-1}) V(Y|\vec{X}_r)$
 - $\varepsilon(r, n) = \mathsf{E}_n \left[(\hat{\beta}_r \beta_r^*)^\top \mathsf{E}(\vec{X}_r \vec{X}_r^\top) (\hat{\beta}_r \beta_r^*) \right]$, where $\hat{\beta}_r$ (or β_r^*) is the training sample (or population) OLS coefficient.

$$\varepsilon(r,n) = \frac{A(r) + \tau^2}{n - r - 2} \left(\frac{n - 2}{n} + r \right)$$

Double Descent without over-fitting



(a) A(p) does not vary with p, $\gamma = p/n$ (Hastie et al. (2019))

Double Descent without over-fitting



Double Descent without over-fitting



Multiple Descents



Regression Tree with Potentially Infinite Depth

• Assume both target and training populations are the same, satisfying

$$X_1, X_2, \ldots \overset{i.i.d.}{\sim}$$
Bernoulli $(1/2), \quad V(Y) < \infty,$

and that dependence of Y on $\{X_1, X_2, \ldots\}$ is arbitrary.

Regression Tree with Potentially Infinite Depth

• Assume both target and training populations are the same, satisfying

$$X_1, X_2, \ldots \stackrel{i.i.d.}{\sim}$$
Bernoulli $(1/2), \quad V(Y) < \infty,$

and that dependence of Y on $\{X_1, X_2, \ldots\}$ is arbitrary.

• Prediction for a unit with covariates value \vec{x}_{∞} :

$$\hat{\mu}_{r}(\vec{x}_{r}) = \begin{cases} \{n(\vec{x}_{r})\}^{-1} \sum_{i:\vec{x}_{ir}=\vec{x}_{r}} Y_{i}, & \text{if } n(\vec{x}_{r}) > 0, \\ \{n(\vec{x}_{k})\}^{-1} \sum_{i:\vec{x}_{ik}=\vec{x}_{k}} Y_{i}, & \text{if } n(\vec{x}_{k}) > 0 \text{ and } n(\vec{x}_{k+1}) = 0, \\ & \text{for some } 0 \le k < r. \end{cases}$$

Regression Tree with Potentially Infinite Depth

• Assume both target and training populations are the same, satisfying

$$X_1, X_2, \ldots \stackrel{i.i.d.}{\sim}$$
Bernoulli $(1/2), \quad V(Y) < \infty,$

and that dependence of Y on $\{X_1, X_2, \ldots\}$ is arbitrary.

• Prediction for a unit with covariates value \vec{x}_{∞} :

$$\hat{\mu}_{r}(\vec{x}_{r}) = \begin{cases} \{n(\vec{x}_{r})\}^{-1} \sum_{i:\vec{x}_{ir}=\vec{x}_{r}} Y_{i}, & \text{if } n(\vec{x}_{r}) > 0, \\ \{n(\vec{x}_{k})\}^{-1} \sum_{i:\vec{x}_{ik}=\vec{x}_{k}} Y_{i}, & \text{if } n(\vec{x}_{k}) > 0 \text{ and } n(\vec{x}_{k+1}) = 0, \\ & \text{for some } 0 \le k < r. \end{cases}$$

• The above is the "highest resolution imputation", given the ordering of the covariates (which means there are other methods.)

Decomposition of the prediction error

• Model error, which is merely the intrinsic error:

$$\tau^2 \equiv \mathsf{E}[\operatorname{Var}(Y \mid \vec{X}_{\infty})]$$

Decomposition of the prediction error

• Model error, which is merely the intrinsic error:

$$au^2 \equiv \mathsf{E}[\operatorname{Var}(Y \mid \vec{X}_{\infty})]$$

• Approximation error:

$$A(r) = \sum_{k=r+1}^{\infty} [\mathsf{E}\{\operatorname{Var}(Y \mid \vec{\boldsymbol{X}}_{k-1})\} - \mathsf{E}\{\operatorname{Var}(Y \mid \vec{\boldsymbol{X}}_{k})\}]$$

Decomposition of the prediction error

• Model error, which is merely the intrinsic error:

$$au^2 \equiv \mathsf{E}[\operatorname{Var}(Y \mid \vec{X}_{\infty})]$$

• Approximation error:

$$A(r) = \sum_{k=r+1}^{\infty} [\mathsf{E}\{\operatorname{Var}(Y \mid \vec{\boldsymbol{X}}_{k-1})\} - \mathsf{E}\{\operatorname{Var}(Y \mid \vec{\boldsymbol{X}}_{k})\}]$$

• Estimation error: $\varepsilon(r, n) = E_n E[\{\hat{\mu}_r(\vec{X}_r) - E(Y \mid \vec{X}_r)\}^2]$

$$\varepsilon(r,n) = \left\{ A(r) + \tau^2 \right\} \cdot \mathsf{E}_n \left[\frac{\mathbbm{1}(n(\vec{1}_r) > 0)}{n(\vec{1}_r)} \right] \\ + \sum_{k=0}^{r-1} \left\{ A(k) + \tau^2 \right\} \cdot \mathsf{E}_n \left[\frac{\mathbbm{1}(n(\vec{1}_k) > 0, n(\vec{1}_{k+1}) = 0)}{n(\vec{1}_k)} \right] \\ + \sum_{k=0}^{r-1} \left\{ A(k) - A(r) \right\} \cdot \mathsf{E}_n \left[\mathbbm{1}(n(\vec{1}_k) > 0, n(\vec{1}_{k+1}) = 0) \right]$$

Xinran Li and Xiao-Li Meng

Stochastic World $\tau^2 > 0$: Optimal resolution R_n when $\mathcal{F}_r = \sigma\{X_1, \ldots, X_r\}$ and with $n \ i.i.d.$ observations

- Resolution Loss: $A(r) = \sum_{i=r}^{\infty} E[(\mu_{i+1} \mu_i)^2 | \mathcal{F}_r]$
- Estimation loss: $\varepsilon(r, n) = \mathsf{E}(\mu_r \hat{\mu}_r)^2$

Table: Optimal R_n and Minimal Loss $L_n \equiv \mathsf{PE}_n - \tau^2$ (as $n \to \infty$)

$\varepsilon(r, n)$ A(r)	$e^{-\xi r}$	$r^{-\xi}$	$\log^{-\xi}(r)$
r^{α}/n	$R_n = c_n \log n$	$c_n n^{1/(\alpha+\xi)}$	$\frac{c_n n^{1/\alpha}}{\log^{\xi/\alpha}(n)}$
(a ≥ 1)	$L_n \propto \log^lpha(n)/n$	$n^{-\xi/(\alpha+\xi)}$	$\log^{-\xi}(n)$
$\frac{\alpha^r}{n}$	$R_n = \frac{\log n + \log c_n}{\xi + \log \alpha}$	$c_n \log(n)$	$c_n \log(n)$
(<i>α</i> > 1)	$L_n \propto n^{-\xi/(\xi+\loglpha)}$	$\log^{-\xi}(n)$	$\{\log \log(n)\}^{-\xi}$

• $c_n = O(1)$ but satisfies different constraints for different settings.

Deterministic World $\tau^2 = 0$, Optimal resolution R_n

$\epsilon(r, n)$ A(r)	$e^{-\xi r}$	$r^{-\xi}$	$\log^{-\xi}(r)$
linear regression	$R_n = n - c_n$ $L_n \propto n e^{-\xi n}$	c _n n n ^{−ξ}	$c_n n^k$, $k \in (0,1]$ $\log^{-\xi}(n)$
regression	$R_n \begin{cases} \gg \text{ or } = c_n \log(n) \\ = c_n \log(n) \\ = c_n \log(n) \end{cases}$	$c_n \log(n)$	$c_n \log(n)$
tree	$L_n \propto egin{cases} n^{-1}, & \xi > \log(2) \ \log(n)/n, & \xi = \log(2) \ n^{-\xi/\log(2)}, & \xi < \log(2) \end{cases}$	$\log^{-\xi}(n)$	$\{\log \log(n)\}^{-\xi}$

Table: Optimal R_n and Minimal Loss $L_n \equiv \mathsf{PE}_n$ (as $n \to \infty$)

• $c_n = O(1)$ but satisfies different constraints for different settings.

Suppose a second ordering {k₀, k₁,..., k_r} includes the *first* r − M_r predictors of the first ordering; i.e., M_r is the minimum integer s.t. {0, 1, 2, ..., r − M_r} ⊂ {k₀,..., k_r}

Suppose a second ordering {k₀, k₁,..., k_r} includes the first r − M_r predictors of the first ordering; i.e., M_r is the minimum integer s.t. {0, 1, 2, ..., r − M_r} ⊂ {k₀,..., k_r}

Then the optimal rate under the second ordering is no worse than that under the first ordering if

• Exponential decaying: $\limsup M_r \leq C$

Suppose a second ordering {k₀, k₁,..., k_r} includes the *first* r − M_r predictors of the first ordering; i.e., M_r is the minimum integer s.t. {0, 1, 2, ..., r − M_r} ⊂ {k₀,..., k_r}

Then the optimal rate under the second ordering is no worse than that under the first ordering if

- Exponential decaying: $\limsup M_r \leq C$
- Polynomial decaying: $\limsup M_r/r < 1$

Suppose a second ordering {k₀, k₁,..., k_r} includes the *first* r − M_r predictors of the first ordering; i.e., M_r is the minimum integer s.t. {0, 1, 2, ..., r − M_r} ⊂ {k₀,..., k_r}

Then the optimal rate under the second ordering is no worse than that under the first ordering if

- Exponential decaying: $\limsup M_r \leq C$
- Polynomial decaying: $\limsup M_r/r < 1$
- Logarithm decaying: $M_r = r r^{1/a_r}$, where $a_r = O(1)$.

Suppose a second ordering {k₀, k₁,..., k_r} includes the *first* r − M_r predictors of the first ordering; i.e., M_r is the minimum integer s.t. {0, 1, 2, ..., r − M_r} ⊂ {k₀,..., k_r}

Then the optimal rate under the second ordering is no worse than that under the first ordering if

- Exponential decaying: $\limsup M_r \leq C$
- Polynomial decaying: $\limsup M_r/r < 1$
- Logarithm decaying: $M_r = r r^{1/a_r}$, where $a_r = O(1)$.

Suppose a second ordering {k₀, k₁,..., k_r} includes the *first* r − M_r predictors of the first ordering; i.e., M_r is the minimum integer s.t. {0, 1, 2, ..., r − M_r} ⊂ {k₀,..., k_r}

Then the optimal rate under the second ordering is no worse than that under the first ordering if

- Exponential decaying: $\limsup M_r \leq C$
- Polynomial decaying: $\limsup M_r/r < 1$
- Logarithm decaying: $M_r = r r^{1/a_r}$, where $a_r = O(1)$.
- The slower A(r) decays, the less restriction on $\{M_r\}$.

Numerical Illustrations Under Linear Models



Э ∃ >

Numerical Illustrations Under Linear Models



Performance of estimators when $\tau^2 > 0$

Table: Estimated resolution and corresponding prediction error when n = 50.

Type / r_{opt} / $PE_n(r_{opt})$	Method	Ŕ	95% QR	$PE_n(\hat{R})/PE_n(r_{\mathrm{opt}})$	95% QR
Exponential	Oracle	_	_	1.00	[0.92, 1.18]
$r_{\rm opt} = 4$	CV	6	[2, 20]	1.34	[0.92, 1.72]
$PE_n(r_{opt}) = 0.5767$	UE	7	[2, 47]	1.93	[0.92, 9.23]
	IC	47	[46, 47]	19.96	[3.86, 81.00]
Polynomial	Oracle	_	_	1.00	[0.87, 1.24]
$r_{\rm opt}=7$	CV	10	[2, 44]	1.37	[0.91, 5.17]
$PE_n(r_{\rm opt}) = 0.7463$	UE	11	[2, 47]	1.62	[0.90, 8.18]
	IC	47	[45, 47]	14.77	[2.93, 57.54]
Logarithmic	Oracle	_	_	1.00	[0.89, 1.21]
$r_{\rm opt}=6$	CV	9	[2, 41]	1.38	[0.91, 3.97]
$PE_n(r_{opt}) = 0.9714$	UE	10	[2, 46]	1.87	[0.91, 8.04]
-	IC	47	[46, 47]	14.60	[3.52, 58.19]

Xinran Li and Xiao-Li Meng

Performance of estimators when $\tau^2 = 0$

Table: Estimated resolution and corresponding prediction error when n = 50.

Type / r_{opt} / $PE_n(r_{opt})$	Method	Ŕ	95% QR	$PE_n(\hat{R})/PE_n(r_{opt})$	95% QR
Exponential	Oracle	_	-	1.08	[0.12, 4.38]
$r_{\rm opt} = 47$	CV	46	[44, 47]	3.57	[0.33, 15.54]
$PE_n(r_{opt}) = 1.90 \times 10^{-19}$	UE	46	[45, 47]	2.30	[0.32, 12.34]
	IC	47	[47, 47]	1.40	[0.29, 6.68]
Polynomial	Oracle	_	_	0.99	[0.71, 1.48]
$r_{\rm opt}=23$	CV	27	[13, 47]	1.31	[0.73, 4.02]
$PE_n(r_{opt}) = 0.0816$	UE	28	[13, 47]	1.44	[0.74, 4.68]
	IC	47	[46, 47]	5.69	[1.28, 23.97]
Logarithmic	Oracle	_	-	1.00	[0.83, 1.34]
$r_{ m opt} = 12$	CV	15	[4, 44]	1.38	[0.87, 5.39]
$PE_n(r_{\mathrm{opt}}) = 0.357$	UE	16	[5, 47]	1.70	[0.87, 7.72]
	IC	47	[46, 47]	10.99	$[2.27, \ 41.35]$

Big data for Individualized Inference/Prediction

• Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.

Big data for Individualized Inference/Prediction

- Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.
- The concept of bias is more critical than (pure) variance, which exists essentially only at the infinite resolution level.

Big data for Individualized Inference/Prediction

- Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.
- The concept of bias is more critical than (pure) variance, which exists essentially only at the infinite resolution level.
- In the stochastic world, a central task is to determine the **appropriate level of approximation via the bias-variance trade-off.**

Big data for Individualized Inference/Prediction

- Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.
- The concept of bias is more critical than (pure) variance, which exists essentially only at the infinite resolution level.
- In the stochastic world, a central task is to determine the **appropriate level of approximation via the bias-variance trade-off.**
- In the deterministic world, we can put all our eggs in the basket of bias, when we have great sparsity.
For those who just woke up ...

Big data for Individualized Inference/Prediction

- Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.
- The concept of bias is more critical than (pure) variance, which exists essentially only at the infinite resolution level.
- In the stochastic world, a central task is to determine the **appropriate level of approximation via the bias-variance trade-off.**
- In the deterministic world, we can put all our eggs in the basket of bias, when we have great sparsity.

For those who just woke up ...

Big data for Individualized Inference/Prediction

- Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.
- The concept of bias is more critical than (pure) variance, which exists essentially only at the infinite resolution level.
- In the stochastic world, a central task is to determine the **appropriate level of approximation via the bias-variance trade-off.**
- In the deterministic world, we can put all our eggs in the basket of bias, when we have great sparsity.

Resolution is a fundamental concept for scientific inference

• Low resolutions render operational meaning to (frequentest) probability and scientific evaluations

For those who just woke up ...

Big data for Individualized Inference/Prediction

- Statistics becomes an approximation scheme: **approximating individuals by proxy populations**, not the other way around.
- The concept of bias is more critical than (pure) variance, which exists essentially only at the infinite resolution level.
- In the stochastic world, a central task is to determine the **appropriate level of approximation via the bias-variance trade-off.**
- In the deterministic world, we can put all our eggs in the basket of bias, when we have great sparsity.

Resolution is a fundamental concept for scientific inference

- Low resolutions render operational meaning to (frequentest) probability and scientific evaluations
- Data, inference, decisions, evaluations may all have different resolutions (Meng, 2021, "Part II")

Xinran Li and Xiao-Li Meng

25

3

メロト メポト メモト メモト

Searching for optimal R_n in the Stochastic World: $\tau^2 > 0$

- $\frac{r_n}{n} = o(1)$ is necessary for $\varepsilon(r_n, n) = o(1)$, under which $\varepsilon(r_n, n) \asymp \frac{r_n}{n}$.
- General Result: Let R_n be a rate-optimal resolution, and $L_n = A(R_n) + \varepsilon(R_n, n)$ be the minimal prediction error (excluding τ^2). Assume polynomial estimation error, that is, $\varepsilon(r, n) \simeq r^{\alpha}/n$.
 - (i) Hard Thresholding: A(r) = 0 for $r \ge r_0$, and A(r) > 0 for $r < r_0$. Then $\overline{R_n \asymp Constant}$ with the constraint that $\liminf_{n \to \infty} R_n \ge r_0$; and $L_n \asymp n^{-1}$.
 - (ii) Exponential $A(r) \simeq e^{-\xi r}$. Then $R_n = a_n \log(n)$ with a_n satisfying $\overline{a_n \simeq 1}$ and $n^{1-\xi a_n} \log^{-\alpha}(n) = O(1)$; and $L_n \simeq n^{-1} \log^{\alpha}(n)$.
 - (iii) Polynomial $A(r) \simeq r^{-\xi}$. Then $R_n \simeq n^{1/(\alpha+\xi)}$; and $L_n \simeq n^{-\xi/(\alpha+\xi)}$.
 - (iv) Logarithmic $A(r) \simeq \log^{-\xi}(r)$. Then $R_n = a_n n^{1/\alpha} \log^{-\xi/\alpha}(n)$ with a_n satisfying $a_n = O(1)$ and $\limsup_{n \to \infty} \left[\log^{-1}(n) \log(a_n^{-1}) \right] < \alpha^{-1}$; and $L_n \simeq \log^{-\xi}(n)$.

26

ヘロト ヘ戸ト ヘヨト ヘヨト

Search for Optimal R_n in the Deterministic World: $\tau^2 = 0$

- The prediction error simplifies to $A(r) \cdot \frac{(n+1)(n-2)}{n(n-r-2)}$. Specific Results (for normal model)
 - (i) Hard Thresholding: A(r) = 0 for $r \ge r_0$, and A(r) > 0 for $r < r_0$. The optimal resolution is any (sequence) R_n such that $\liminf_{n \to \infty} R_n \ge r_0$ and

$$R_n \leq n-3$$
; and $L_n = 0$.

- (ii) Exponential $A(r) \simeq e^{-\xi r}$. $R_n = n O(1)$ with $R_n \le n 3$; and $\overline{L_n \simeq n e^{-\xi n}}$.
- (iii) Polynomial $A(r) \simeq r^{-\xi}$. $R_n = a_n n$ with a_n satisfying $a_n \simeq 1$ and $\lim \sup a_n < 1$; and $L_n \simeq n^{-\xi}$.
- (iv) Logarithmic $A(r) \simeq \log^{-\xi}(r)$. Optimal resolution is any (sequence) R_n such that $\limsup R_n/n < 1$, $\liminf \frac{\log R_n}{\log n} > 0$; and $L_n \simeq \log^{-\xi}(n)$.

Search for Optimal R_n in the Stochastic World: $\tau^2 > 0$

- $\frac{2^r}{n} = o(1)$ is necessary for $\varepsilon(r, n) = o(1)$, under which $\varepsilon(r, n) \asymp \frac{2^r}{n}$.
- General Result: Assume exponential estimation error: $\varepsilon(r, n) \simeq \frac{\alpha'}{n}$.
 - (i) Hard Thresholding: A(r) = 0 for $r \ge r_0$, and A(r) > 0 for $r < r_0$. Then $R_n \simeq Constant$ with the constraint that $\liminf_{n \to \infty} R_n \ge r_0$, and $L_n \simeq n^{-1}$.
 - (ii) Exponential $A(r) \simeq e^{-\xi r}$. Then $R_n = \frac{\log(n) + \log(a_n)}{\log(\alpha) + \xi}$ with $a_n \simeq 1$; and $L_n \simeq n^{-\xi/\{\log(\alpha) + \xi\}}$.
 - (iii) Polynomial $A(r) \simeq r^{-\xi}$. Then $R_n = a_n \log(n)$ with a_n satisfying $a_n \simeq 1$ and $n^{a_n \log(\alpha) - 1} \log^{\xi}(n) = O(1)$; and $L_n \simeq \log^{-\xi}(n)$.
 - (iv) Logarithmic $A(r) \approx \log^{-\xi}(r)$. Then $R_n = a_n \log(n)$ with a_n satisfying

$$\liminf_{n\to\infty} \frac{\log(a_n)}{\log\log(n)} > -1, \quad \text{and} \quad \frac{\{\log\log(n)\}^{\xi}}{n^{1-a_n\log(\alpha)}} = O(1);$$

and $L_n \asymp \{\log \log(n)\}^{-\xi}$.

Search optimal R_n in the deterministic World: $\sigma^2 = 0$

Let
$$L_n = A(R_n) + \varepsilon(R_n, n) \le A(R_n) + \overline{\varepsilon}(R_n, n) \equiv \overline{L}_n$$
. Specific Results:
(i) Hard Thresholding: $A(r) = 0$ for $r \ge r_0$, and $A(r) > 0$ for $r < r_0$.
Then R_n satisfies that $\liminf_{n \to \infty} R_n \ge r_0$; and $L_n \asymp (1 - 2^{-r_0})^n$.
(ii) Exponential $A(r) \asymp e^{-\xi r}$.
(a) If $e^{-\xi} < 1/2$, then \overline{R}_n satisfies $ne^{-\xi R_n} = O(1)$; and $\overline{L}_n \asymp n^{-1}$.
(b) If $e^{-\xi} = 1/2$, then $\overline{R}_n = a_n \log(n)$ with a_n satisfying $a_n \asymp 1$ and $n^{1-a_n \log(2)}/\log(n) = O(1)$; and $\overline{L}_n \asymp n^{-1} \log(n)$.
(c) If $e^{-\xi} > 1/2$, then $\overline{R}_n = a_n \log(n)$ with a_n satisfying $n^{a_n-1/\log(2)} \asymp 1$;
and $\overline{L}_n \asymp n^{-\xi/\log(2)}$.
(iii) Polynomial $A(r) \asymp r^{-\xi}$. Then $\overline{R}_n = a_n \log(n)$ with a_n satisfying $a_n \asymp 1$ and $n^{a_n \log(2)-1} = O(1)$; and $\overline{L}_n \asymp \log^{-\xi}(n)$.

(iv) Logarithmic $A(r) \approx \log^{-\xi}(r)$. Then $\overline{R}_n = a_n \log(n)$ with a_n satisfying

$$\liminf_{n\to\infty} \frac{\log(a_n)}{\log\log(n)} > -1, \quad \text{and} \quad n^{a_n\log(2)-1} = O(1);$$

and $\overline{L}_n \simeq \{\log \log(n)\}^{-\xi}$.

An asymptotic lower bound for the estimation error $\varepsilon(r, n)$ has the following form.

(i) Hard Thresholding:
$$A(r) = 0$$
 for $r \ge r_0$, and $A(r) > 0$ for $r < r_0$.
Then $A(r_n) + \varepsilon(r, n) \gtrsim (1 - 2^{-r_0})^n$.

(ii) Exponential $A(r) \simeq e^{-\xi r}$. Then $A(r_n) + \varepsilon(r, n) \gtrsim n^{-\xi/\log(2)}$.

(iii) Polynomial
$$A(r) \simeq r^{-\xi}$$
. Then $A(r_n) + \varepsilon(r, n) \gtrsim \log^{-\xi}(n)$.

(iv) Logarithmic $A(r) \asymp \log^{-\xi}(r)$. Then $A(r_n) + \varepsilon(r, n) \gtrsim \{\log \log(n)\}^{-\xi}$.