## Advances in Survey Estimation with Imperfectly Matched Auxiliary Data

# Jay Breidt Colorado State University

#### 6th Italian Conference on Survey Methodology June 7, 2019

Joint work with Chien-Min Huang and Jean Opsomer Colorado State University and Westat

- About 450 charter boats and 15,000 boat trips along the Atlantic Coast of South Carolina each year
- How many black sea bass were caught in 2018?

 $U = \{1, 2, \dots, N\}$ 

- = {all SC charter fishing boat trips in 2018}
- -number of black sea bass caught on kth trip:  $y_k$
- -total black sea bass caught:  $T = \sum_{k \in U} y_k$
- $\bullet$  Infeasible to obtain data on all  $N\simeq 15,000$  boat trips: instead, use a probability sample  $s\subset U$

## Two sources of information on the charter boat fishery 3

### Sample with angler interviews: Monthly logbook records:



	al (Diacas Dais		500	THCAR	JLINA CHARL	EKB	Date	001	·	D	Revised	4-2012
vess	ei (Piease Prir	n):					Date: _			Pern	iit iso.:	
Nun	tber of Angler	s:	Tr	ip Start Ti	me: Ac	tual l	Hours Fishe	d: _		Locatio	DI: Example: 32-	78-C1 (see
Trip	Start ation:			Arti —— Ree	ficial f Name:				Targe Sneci	t es:	Example: 52	o er (se
Log	alas 🗖 Estur	rina		Mathody		Cast	/ El.	Wate	e Donti	(Plea	se specifiy)	6
LOCI	are: $\Box = 10 - 3$	miles		Methou:		Dive	. / Fiy	wate	r Depu	i: Shan	owest:	0
Offshore     Offshore							Dive Gig Deepest:					I
		1010			AGENCY USE	ONLY	2					
	MAIL OR FA	X RE	PORT B	Y	Yr Mo	Da	ay Permit	#		Location	Locale	Ang# N
601	THE 10 <sup>th</sup> OF 7	THE M	ONTH 1	:0: 								
Bo	x 12559 Charle	stausu	C 29423	a, P.O.	Target Sp.	Hrs	. Reef	1	rip Start		Shallowest	Deepest
FAX	(843) 953-936	2 Phone	e: (843) 9	53-9313								
_		#	Lbs	# Release	d # Released					Lbs	# Released	# Relea
	Species	Kept	Kept	Alive	Dead		Species		# Kept	Kept	Alive	Dea
1050	Dolphin				_	1423	Gag					
4/10	wahoo					1424	Scamp					
4655	Yellowin Tuna					1414	Snowy Group	er				
4658	Blackfin Tuna					1416	Red Grouper					
3026	Sailfish					1410	Other Group	er				
2177	White Marlin					2202	(Specify)					
2179	Blue Marlin				_	3302	Red Porgy (P	nks)				
1940	King Mackeral				_	3295	Other Porgies	5				
3840	Spanish Mackeral						(Specify)					
4653	Little Tunny					3764	Red Snapper					
0330	Bonita				_	3765	Vermillion Sna	pper				
4654	Skip Jack					3360	0 Black Sea Bass					
0180	Barracuda					3314	Spottail Pinte	sh				
3810	Spadefish				_	1441	White Grunt					
0030	Amberjack				_	1440	Other Grunts					
0870	Crevalle Jack				_		(Specify)					
0230	Bluefish					4560	Triggerfish					
0570	Cobia					1082	Red Drum					
4350	Tarpon					1081	Black Drum					
	Other Fish				_	3447	Spotted Seatr	out				
	(Specify)					3446	Weakfish					
						1209	Flounder					
					-	3560	Sneepshead					
_						4410	Ladyrish					
Canto	in'e Notae-		L			2670	Inchore Plate	ah				
Captain's Potes:						3518	Sharnnose Sh	ark				
						3495	Blacktin Sho	*				
						3483	Bonnethead S	 hark		-		
						3521	Spiny Dogfist	1				
						3511	Smooth Dogfi	sh			1	
Signature:					3508	Other Sharks						
						(Specify)			1			
							1 · · · · · ·			1		

## Goal: combine logbook database with survey data

- Design-based difference estimators
- Extension to multiple frames
- Require matching of sampled elements to auxiliary records
  - most theory and methods assume matching is done without error
- Some results on estimation under imperfect matching
  - properties of difference estimators
  - simulation results based on South Carolina charter boat fishing

- Draw probability sample  $s \subset U$  via design with known, positive inclusion probabilities  $\Pr[k \in s] = \pi_k > 0$
- $\bullet$  Sample membership indicator  $I_k=1$  if  $k\in s,\ I_k=0$  otherwise

$$\mathsf{E}\left[I_k\right] = \pi_k$$
, averaging over all possible samples

• Since  $E[I_k/\pi_k] = 1$  under repeated sampling, unbiased Horvitz-Thompson estimator of T is

$$\widehat{T} = \sum_{k \in S} \frac{y_k}{\pi_k} = \sum_{k \in U} y_k \frac{I_k}{\pi_k}$$

- ullet Auxiliary data  $oldsymbol{x}_\ell$  for all  $\ell$  in some database  $\mathcal{A}$
- Perfect, known matching from  $\mathcal{A}$  to population U:

$$M_{k\ell} = \begin{cases} 1, & \text{if } \ell \in \mathcal{A} \text{ matches } k \in U, \\ 0, & \text{otherwise} \end{cases}$$

• A "method"  $\mu(\cdot)$  for predicting  $y_k$  from  $oldsymbol{x}_\ell$ :

$$\sum_{\ell \in \mathcal{A}} M_{k\ell} \, \mu(\boldsymbol{x}_\ell) = \widetilde{y}_k \text{ predicts } y_k$$

- for each element k, look up the correct  $oldsymbol{x}_\ell$
- apply  $\mu(\cdot)$ , which does not depend on the sample

### Difference estimator combines sample and auxiliary data 7

 $\bullet$  Difference estimator of T is then

$$\begin{split} \widetilde{T} &= \sum_{k \in U} \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_{\ell}) + \sum_{k \in s} \frac{y_k - \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_{\ell})}{\pi_k} \\ &= \sum_{k \in U} \widetilde{y}_k + \sum_{k \in U} (y_k - \widetilde{y}_k) \frac{I_k}{\pi_k} \\ &= \text{(auxiliary-based prediction)} + \text{(bias adjustment)} \end{split}$$

where  $\widetilde{y}_k$  is not random

• Expectation is

$$\mathsf{E}\left[\widetilde{T}\right] = \sum_{k \in U} \widetilde{y}_k + \sum_{k \in U} \left(y_k - \widetilde{y}_k\right) \mathsf{E}\left[\frac{I_k}{\pi_k}\right] = T$$

$$\operatorname{Var}\left(\sum_{k\in U}\widetilde{y}_{k} + \sum_{k\in U}(y_{k} - \widetilde{y}_{k})\frac{I_{k}}{\pi_{k}}\right)$$
$$= \sum_{j,k\in U}\Delta_{jk}\frac{(y_{j} - \widetilde{y}_{j})(y_{k} - \widetilde{y}_{k})}{\pi_{j}}\frac{(y_{k} - \widetilde{y}_{k})}{\pi_{k}}$$

• Compare to Horvitz-Thompson estimator:

$$\operatorname{Var}\left(\sum_{k\in U} y_k \frac{I_k}{\pi_k}\right) = \sum_{j,k\in U} \Delta_{jk} \frac{y_j y_k}{\pi_j \pi_k}$$

- $\bullet$  Difference estimator is exactly unbiased, regardless of the quality of the method  $\mu(\cdot)$
- Has smaller variance than HT provided "residuals"

$$y_k - \widetilde{y}_k$$

have smaller variation than "raw values"  $y_k$ 

$$-(\text{If } M_{k\ell} \equiv 0, \text{ we get back HT})$$

- Have an exactly unbiased variance estimator
- Above results assume (1) one frame covers the universe and (2) matching is perfect

• Assume that the universe U is completely covered by disjoint "overlap domains":

 $U = \left\{ \cup_{g \in G_1} U_g \right\} \cup \left\{ \cup_{g \in G_2} U_g \right\} \cup \left\{ \cup_{g \in G_3} U_g \right\}$ 

- $\bullet$  If  $g \in G_1,$  overlap domain  $U_g$  is covered by one or more frames, but not the database
- $\bullet$  If  $g \in G_2$ , overlap domain  $U_g$  is covered by one or more frames and the database
- $\bullet$  If  $g \in G_3,$  overlap domain  $U_g$  is covered only by the database



		In Auxiliary Database?						
		Νο	Yes					
In Sampling Frame(s)?	No		<ul> <li>G<sub>3</sub></li> <li>Synthetic predictor</li> <li>Biased</li> <li>Zero sampling variance</li> </ul>					
	Yes	<ul> <li>G<sub>1</sub></li> <li>Mecatti estimator</li> <li>Unbiased</li> <li>Potentially large variance</li> </ul>	<ul> <li>G<sub>2</sub></li> <li>Difference estimator</li> <li>Unbiased</li> <li>Small variance if auxiliary information is good</li> </ul>					

- From frame f, draw a sample  $s_{fg}$  to represent  $U_g$
- Compute Horvitz-Thompson estimator

$$\widehat{T}_{fg} = \sum_{k \in s_{fg}} \frac{y_k}{\pi_k^{(f)}}, \quad \text{where } \mathsf{E}\left[\widehat{T}_{fg}\right] = T_g$$

• Define the coverage indicator

 $F_{fg} = \begin{cases} 1, & \text{if overlap domain } g \text{ is covered by frame } f \\ 0, & \text{otherwise} \end{cases}$ 

• Adjust for multiplicity by constructing weights

$$\psi_{fg} = \frac{F_{fg}}{\left(\sum_{f} F_{fg}\right)}$$

( $\psi_{fg} = 1$  if domain covered by only one frame; 1/2 if two frames, etc.)

• Unbiased Mecatti/multiplicity estimator for  $\sum_{g \in G_1} T_g$  is



• Multiplicity-adjusted difference estimator for  $g \in G_2$ :

$$\widetilde{T}_{g}^{*} = \sum_{k \in U_{g}} \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_{\ell}) + \sum_{f=1}^{F} \psi_{fg} \sum_{k \in s_{fg}} \frac{y_{k} - \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_{\ell})}{\pi_{k}^{(f)}}$$
$$= \sum_{k \in U_{g}} \widetilde{y}_{k} + \sum_{f=1}^{F} \psi_{fg} \sum_{k \in U_{g}} (y_{k} - \widetilde{y}_{k}) \frac{I_{k}^{(f)}}{\pi_{k}^{(f)}}$$

• Unbiased difference estimator for  $\sum_{g \in G_2} T_g$  is then

$$\sum_{g \in G_2} \widetilde{T}_g^*$$

- $G_3$  has no sampling frame coverage
- Can only predict with the auxiliary data,

$$\widetilde{T}_g = \sum_{k \in U_g} \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_\ell) = \sum_{k \in U_g} \widetilde{y}_k$$

 $\bullet$  Synthetic predictor for  $\sum_{g\in G_3} T_g$  is then

$$\sum_{g \in G_3} \widetilde{T}_g = \sum_{g \in G_3} \sum_{k \in U_g} \widetilde{y}_k$$

• Zero sampling variance, unknown bias

- Replace  $M_{k\ell} = 0$  or 1 by match metrics  $m_{k\ell} \in [0, 1]$ - known only for sampled k
- Produced by deterministic algorithm
- Could involve formal probabilistic record linkage (Fellegi and Sunter 1969, Winkler 2009) or other methods

- conditional probabilities, likelihood ratios, ...

• Whatever their origin, treat  $m_{k\ell}$  as fixed in what follows

• Under perfect matching, multi-frame estimator is

$$\sum_{g \in G_1} \sum_{f=1}^F \psi_{fg} \widehat{T}_{fg} + \sum_{\ell \in \mathcal{A}} \left( \sum_{g \in G_2 \cup G_3} \sum_{k \in U_g} M_{k\ell} \right) \mu(\boldsymbol{x}_{\ell}) \\ + \sum_{g \in G_2} \sum_{f=1}^F \psi_{fg} \sum_{k \in s_{fg}} \frac{y_k - \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_{\ell})}{\pi_k^{(f)}}$$

- Under imperfect,  $m_{k\ell}$  is known only for  $k \in s_{fq}$
- $\bullet$  Cannot just substitute  $m_{k\ell}$  for  $M_{k\ell}$  in second term, but ok in third

• Second term under perfect matching is

$$\sum_{\ell \in \mathcal{A}} \left( \sum_{g \in G_2 \cup G_3} \sum_{k \in U_g} M_{k\ell} \right) \mu(\boldsymbol{x}_\ell)$$

- If  $\ell$ th record matches some element in  $\bigcup_{g \in G_2 \cup G_3} U_g$ , then (parenthetical term) = 1
- Under imperfect matching, estimate parenthetical term as equal to 1
  - (or construct a complicated, and biased, estimator)

• Analogue of perfect-match multi-frame estimator

$$\sum_{g \in G_1} \sum_{f=1}^F \psi_{fg} \widehat{T}_{fg} + \sum_{\ell \in \mathcal{A}} \left( \sum_{g \in G_2 \cup G_3} \sum_{k \in U_g} M_{k\ell} \right) \mu(\boldsymbol{x}_{\ell}) \\ + \sum_{g \in G_2} \sum_{f=1}^F \psi_{fg} \sum_{k \in s_{fg}} \frac{y_k - \sum_{\ell \in \mathcal{A}} M_{k\ell} \mu(\boldsymbol{x}_{\ell})}{\pi_k^{(f)}}$$

is then

$$\widetilde{T}_{diff} = \sum_{g \in G_1} \sum_{f=1}^{F} \psi_{fg} \widehat{T}_{fg} + \sum_{\ell \in \mathcal{A}} (1) \mu(\boldsymbol{x}_{\ell}) \\ + \sum_{g \in G_2} \sum_{f=1}^{F} \psi_{fg} \sum_{k \in s_{fg}} \frac{y_k - \sum_{\ell \in \mathcal{A}} m_{k\ell} \mu(\boldsymbol{x}_{\ell})}{\pi_k^{(f)}}$$

• Bias depends on matching and prediction error:

$$\mathsf{E}\left[\widetilde{T}_{diff}\right] - T = -\sum_{g \in G_3} T_g + \sum_{\ell \in \mathcal{A}} \mu(\boldsymbol{x}_{\ell}) - \sum_{\ell \in \mathcal{A}} \left(\sum_{g \in G_2} \sum_{k \in U_g} m_{k\ell}\right) \mu(\boldsymbol{x}_{\ell})$$
  
= -(total uncovered) + (database)-(overlap)

1

• Sufficient conditions for unbiased estimation are

$$G_3 = \emptyset$$
 and  $\sum_{g \in G_2} \sum_{k \in U_g} m_{k\ell} = 1$  for all  $\ell \in \mathcal{A}$ 

 Asymptotic unbiasedness and mean square consistency requires some combination of "not too much" matching error or undercoverage, and "good" prediction of the uncovered population

 $\mathbf{N}$ 

• Variance of the estimator is (setting  $m_{k\ell} \equiv 0$  for  $k \in \bigcup_{g \in G_1} U_g$ ):

$$\begin{split} \sum_{f=1}^{F} \sum_{g \in G_1 \cup G_2} \sum_{g' \in G_1 \cup G_2} \psi_{fg} \psi_{fg'} \sum_{j \in U_g} \sum_{k \in U_{g'}} \Delta_{jk}^{(f)} \frac{d_j}{\pi_j^{(f)}} \frac{d_k}{\pi_k^{(f)}} \\ \text{with } d_j &= \begin{cases} y_j - \sum_{\ell \in \mathcal{A}} M_{j\ell} \mu(\boldsymbol{x}_\ell), & \text{perfect matching} \\ \text{prediction error} \\ y_j - \sum_{\ell \in \mathcal{A}} m_{j\ell} \mu(\boldsymbol{x}_\ell), & \text{imperfect matching} \\ \text{matching and/or prediction error} \end{cases} \\ \bullet \operatorname{Var}\left(\widetilde{T}_{diff}\right) = O\left(\frac{N^2}{\min_f n_f}\right) \text{ and } N^{-1}\widetilde{T}_{diff} \xrightarrow{\mathrm{m.s.}} N^{-1} \mathsf{E}\left[\widetilde{T}_{diff}\right] \end{split}$$

 $\bullet$  Unbiased variance estimation provided all  $\pi^{(f)}_{jk}>0$  in each frame

### Use SC recreational fishery to devise a simulation study 23

- About 450 charter boats and 15,000 boat trips along the Atlantic Coast each year
- Survey data from sampled angler on boat trip on the actual date
  - coverage error: not all sites and times are in-frame
  - lots of sampling error
- Logbook data from captain's report, later that month
  - nonresponse
  - measurement error
- Lots of matching error!

- Perfect match:  $m_{k\ell} = 1$  for  $\ell = \ell_1$  and 0 otherwise
- High-quality match:  $m_{k\bullet} = \sum_{\ell \in \mathcal{A}} m_{k\ell} = 1$

$$m_{k\ell} = \begin{cases} 1/3, & \text{if } \ell = \ell_1, \ \ell = \ell_2 \text{ or } \ell = \ell_3, \\ 0, & \text{otherwise} \end{cases}$$

• Low-quality match:  $m_{k\bullet} = \sum_{\ell \in \mathcal{A}} m_{k\ell} < 1$ 

$$m_{k\ell} = \begin{cases} 1/6, & \text{if } \ell = \ell_1, \ \ell = \ell_2 \text{ or } \ell = \ell_3, \\ 0, & \text{otherwise} \end{cases}$$

• No match:  $m_{k\ell} = 0$  for all  $\ell \in \mathcal{A}$ 

• Match metrics  $\{m_{k\ell}\}_{k\in s,\ell\in\mathcal{A}}$  developed by South Carolina Department of Natural Resources staff

Interview variables	Logbook variables
Date of interview	Date of reported trip
Time of interview	Estimated trip end time
License number of vessel	License number of vessel
Name of vessel given	Name of vessel reporting
Interview site	Reported start site

• Large fraction of unmatched trips and low-quality matches

	No Match	LQ	HQ	Perfect
Empirical	11.0%	52.5%	36.5%	0.0%

- Use real logbook data to create artificial population with |U| = 10,647 boat trips, sorted in space and time
- Use Markov chain to assign (unobservable) states to groups of population boat trips:

	state from Markov chain				
	no match	LQ	HQ	perfect	
size of group of elements:	1	10	5	1	
logbook records created:	0	5	5	1	
metric sum:	0	1/2	1	1	

• If an LQ element is selected, metrics (correctly) indicate it might match one of five records, or none of them

• Set Markov chain parameters to simulate match metrics  $\{m_{k\ell}\}$  and logbook database  $\mathcal{A}$  under two scenarios:

	$ \mathcal{A} $	No Match	LQ	HQ	Perfect
Poor Match	6,836	8.6%	54.4%	31.7%	5.3%
Better Match	9,031	2.3%	23.3%	69.8%	4.7%
Empirical		11.0%	52.5%	36.5%	0.0%

- Population of boat-trips and database of logbook records is then fixed
- Create two incomplete frames, partially overlapping
- Sample repeatedly from this finite population

- Simplifications:
  - $-\operatorname{no}$  "differential matching": quality of  $m_{k\ell}$  does not depend on  $y_k$
  - no measurement error:  $\mu(oldsymbol{x}_\ell) = y_k$  for perfect match
- Draw 1000 repeated samples from simulated population - stratified, two-stage, unequal-probability selection
- Compute  $\widehat{T}_{HT,1}$ ,  $\widehat{T}_{HT,2}$ ,  $\widehat{T}_{Mec}$ ,  $\widetilde{T}_{diff}$  for number of angler trips and several species in each simulated sample
- Assess bias, variance, and MSE for each estimator

#### Even with poor match, difference dominates Mecatti 29



**Angler Trips** 

#### Real Black Sea Bass logbook, HT, and combinations 30

• Frequently targeted and caught; appears regularly in both sources



Black Sea Bass SC 2016





**Black Sea Bass Catch** 

- Auxiliary information is useful even with imperfect matching
  - naive difference estimator improves accuracy and precision of multiplicity estimator
  - variance estimators and confidence intervals (not shown) work well
- Matching across frames or matching across auxiliary databases adds challenges
- Grazie mille!

Contact info: FJay.Breidt@colostate.edu