

# Integrating Reward Maximization and Population Estimation

SEQUENTIAL DECISION-MAKING FOR INTERNAL REVENUE SERVICE AUDIT SELECTION

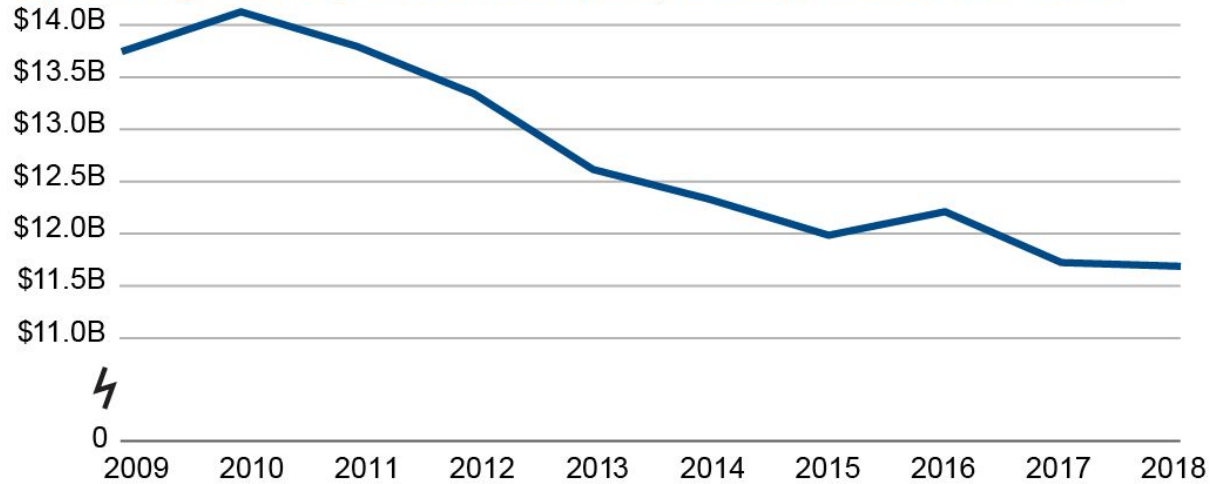
Peter Henderson<sup>1</sup>, Ben Chugg<sup>1</sup>, Brandon Anderson<sup>1,2</sup>, Kristen Altenburger<sup>1</sup>, Alex Turk<sup>2</sup>, John Guyton<sup>2</sup>, Jacob Goldin<sup>1</sup>, Daniel E. Ho<sup>1</sup>

<sup>1</sup>Stanford University <sup>2</sup>IRS RAAS

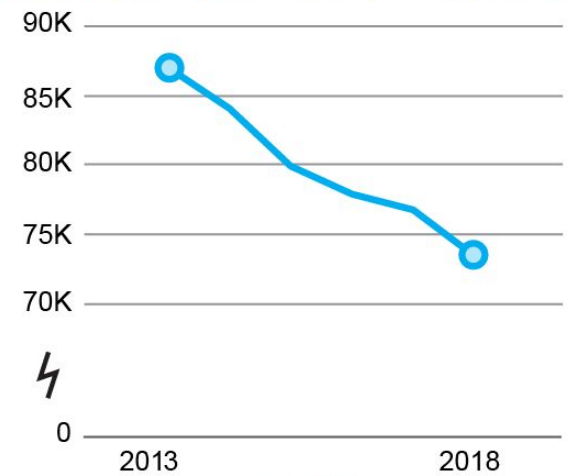
All views and opinions expressed in this presentation are our own and not of any of our co-authors, nor of the Internal Revenue Service or any other company or government entity.

# Institutional Context

## Operating Costs (Constant Dollars), Fiscal Years 2009–2018

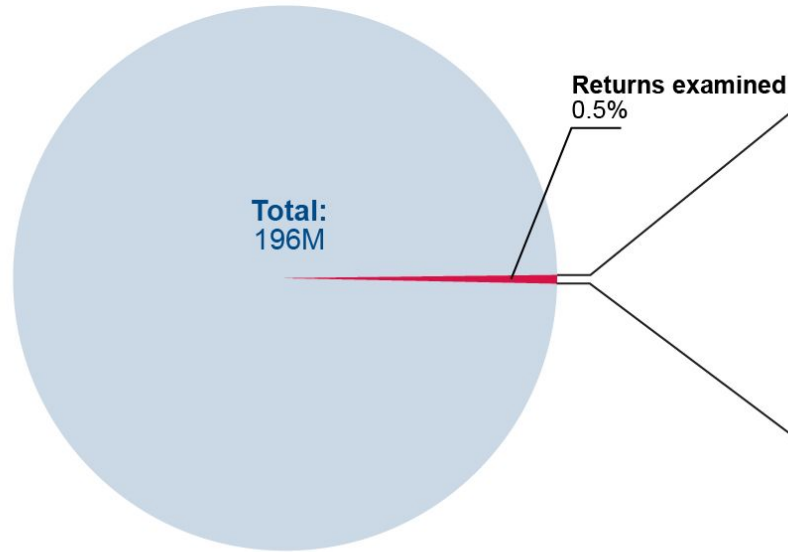


## Full-time Equivalent Positions Realized, Fiscal Years 2013–2018

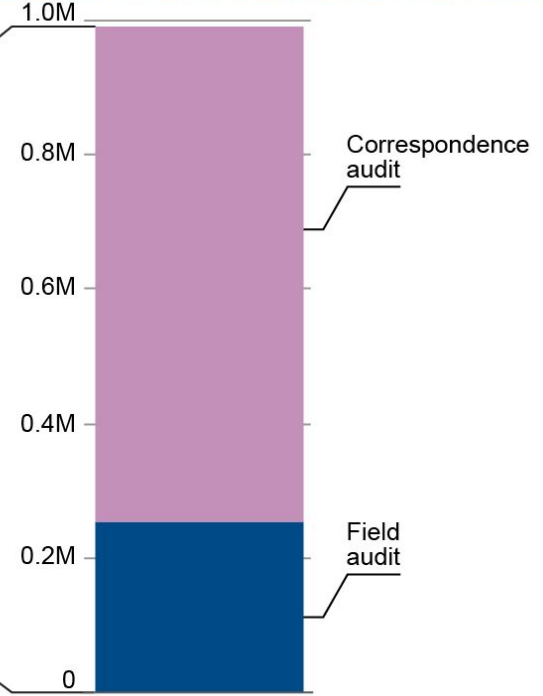


# Institutional Context

Total Returns Filed, Calendar Year 2017, and Percentage Examined, Fiscal Year 2018



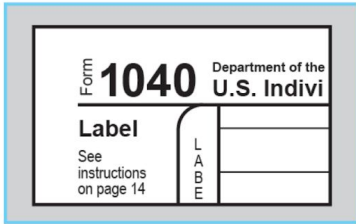
Returns Examined, Fiscal Year 2018



SOURCE: 2018 IRS Data Book Table 9a

# Stylized Program

Identify returns



Form **1040** Department of the U.S. Individual Income Tax Return

**Label**  
See instructions on page 14

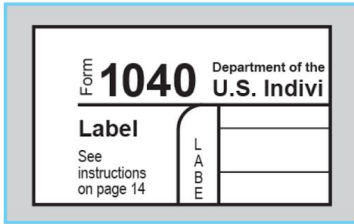
L  
A  
B  
E




Random sample  
(~15k / year, 2006-14)

# Stylized Program

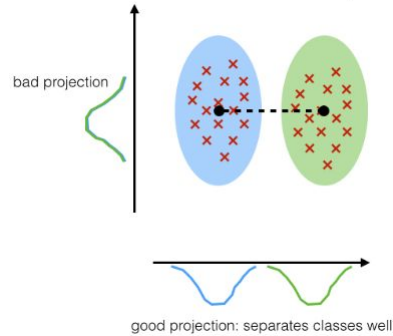
Identify returns



Form 1040 Department of the U.S. Indivi

Label	L	
See instructions on page 14	A	
	B	
	E	

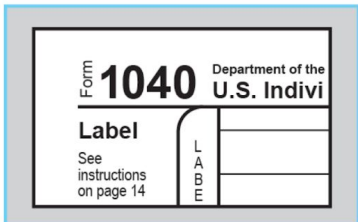
Random sample  
(~15k / year, 2006-14)



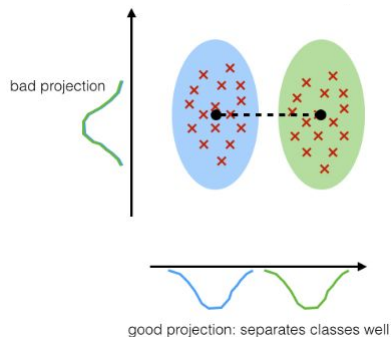
Risk model

# Stylized Program

## Identify returns



Random sample  
(~15k / year, 2006-14)



Risk model



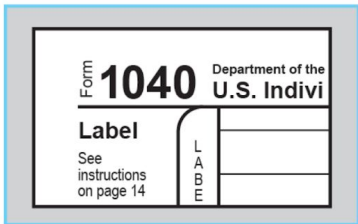
## Audit returns



Risk selected Op Audits  
(>500k / year)

# Stylized Program

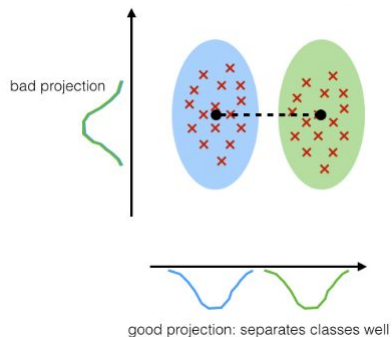
## Identify returns



Random sample  
(~15k / year, 2006-14)



Tax Gap Estimate



Risk model

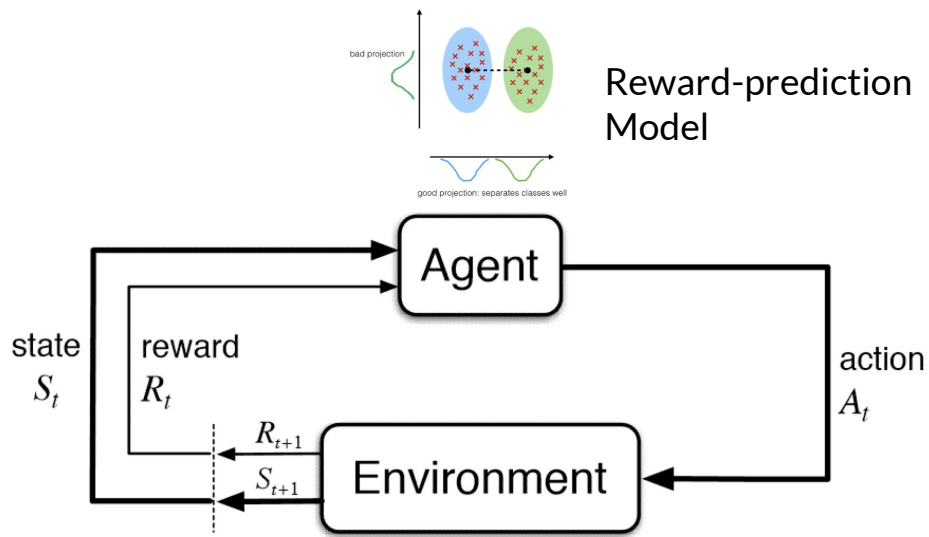


## Audit returns



Risk selected Op Audits  
(>500k / year)

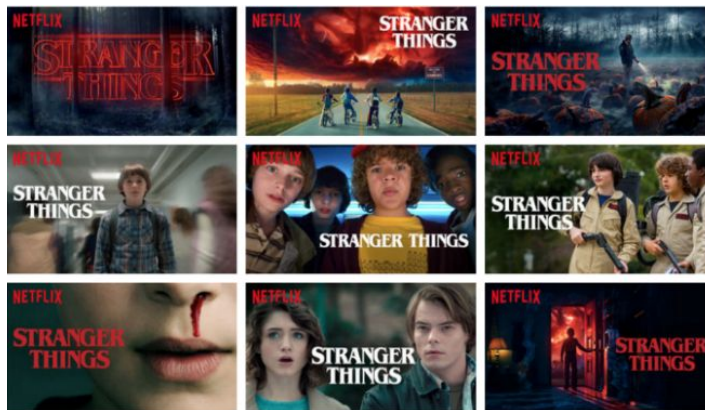
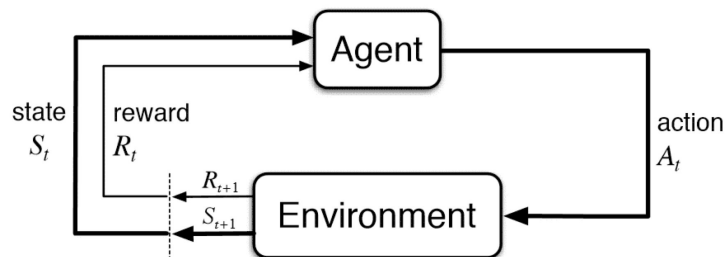
# Sequential Decision-Making (Machine Learning)





# Sequential Decision-Making In the Real World

Example: **NETFLIX**

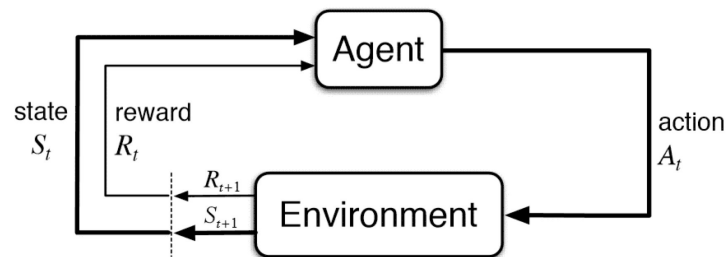


Context	User information on device (environment)
Actions	Set of movie banners to show
Reward	User engagement (click-through, minutes)
Learner	Identify policy to maximize cumulative reward

**Explore** new movies / preferences vs. **Exploit** known preferences

# Sequential Decision-Making In the Real World

Example:  IRS



Identify returns

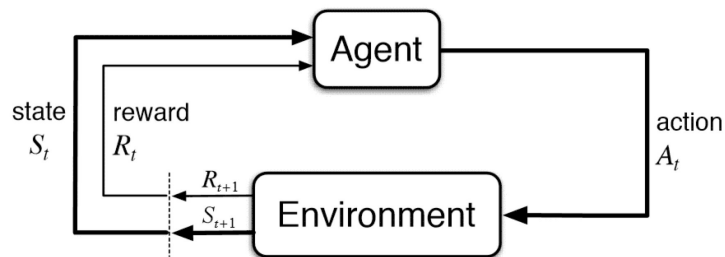
The image shows a tax form label for Form 1040, Department of the U.S. Individual Income Tax Service. The label includes the text "Form 1040 Department of the U.S. Individual Income Tax Service" and "Label See instructions on page 14". The label is divided into sections, with the word "LABEL" written vertically in a column.

Context	Tax return information (taxpayer, stratum, etc.)
Actions	Selecting returns to audit
Reward	Detected Under-reporting
Learner	Identify policy to maximize cumulative reward

**Explore** forms of underreporting vs. **Exploit** known underreporting

# Sequential Decision-Making In the Real World

Example:  IRS



Identify returns

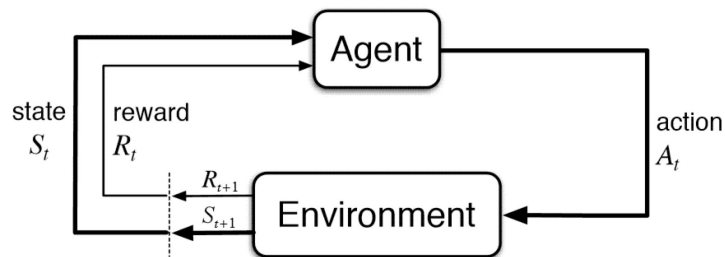
The image shows a tax form label for Form 1040, Department of the U.S. Individual Income Tax. The label includes the text "Form 1040 Department of the U.S. Individual Income Tax" and "Label See instructions on page 14". The label is divided into sections, with the word "LABEL" written vertically in a column.

Context	Tax return information (taxpayer, stratum, etc.)
Actions	Selecting returns to audit
Reward	Under-reporting
Learner	Identify policy to maximize cumulative reward

+ Estimate unbiased population statistics (e.g., tax gap, average misreporting)

# Sequential Decision-Making In the Real World

Example:  IRS



Identify returns

The image shows a tax form label for Form 1040, Department of the U.S. Individual Income Tax. The label includes the text "Form 1040 Department of the U.S. Individual Income Tax" and "Label See instructions on page 14". The label is divided into sections, with the word "LABEL" written vertically in a column.

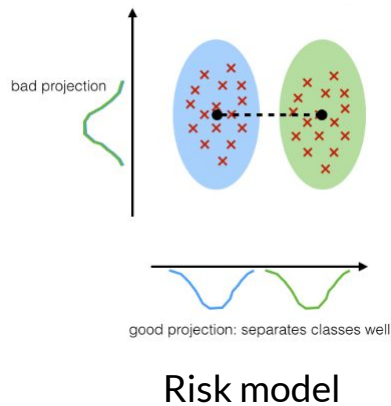
Secondary objective  
not typical of machine  
learning literature



+ Estimate unbiased population statistics (e.g., tax gap, average misreporting)

# Sequential Decision-Making In the Real World

Example:  IRS



**Tempting Solution:** Use a regression-based risk-model to do selection and estimation, with no random sampling.

**Problems:** Sequentially-learned models are known to be biased and there are no theoretically guaranteed ways to *remove* this bias in the low sample regime (yet). (Nie et al., 2018)

Lack of exploration leads to suboptimal feedback loops. (Jiang et al., 2019)

# Optimize-and-Estimate Structured Bandits

Machine Learning Literature  
on Sequential Decision-Making  
(e.g., bandit algorithms that  
optimize for reward only)

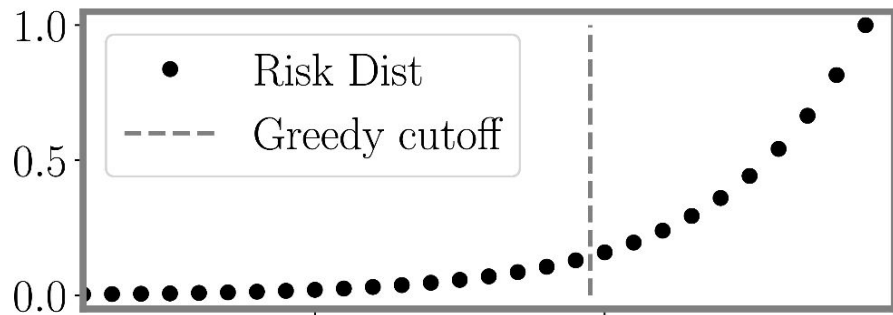
+

Sampling Literature  
(unbiased estimation of  
population statistics)

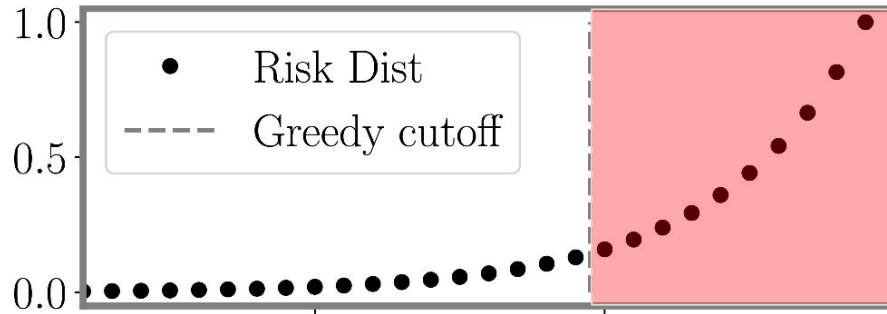
=

Optimize-and-estimate Structured Bandits

# Adaptive Bin Sampling



# Adaptive Bin Sampling

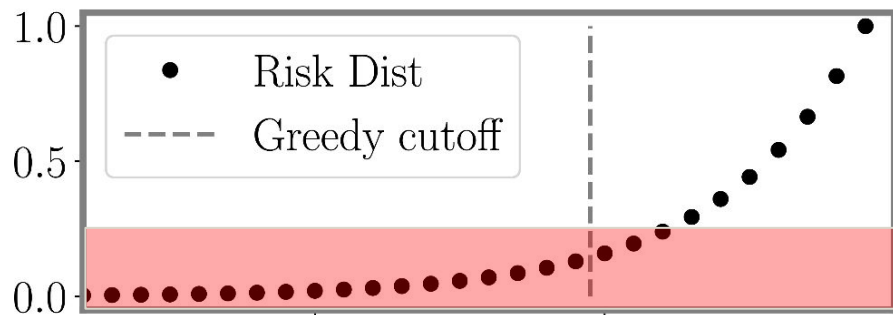


Greedy selection  
(e.g., stylized version of Op audits)

If only use this:  
biased model, biased estimate



# Adaptive Bin Sampling

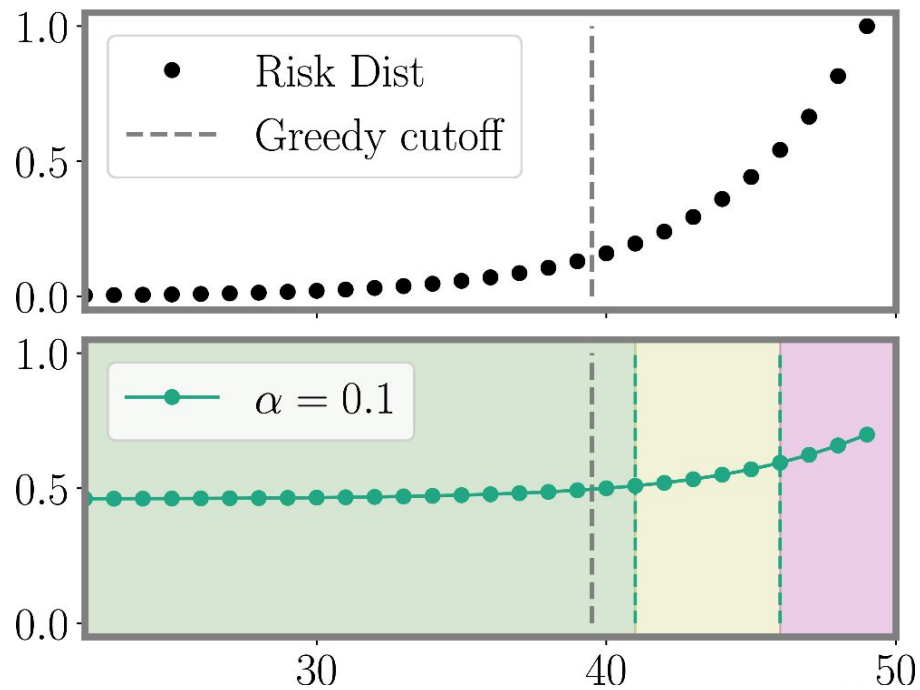


Random selection

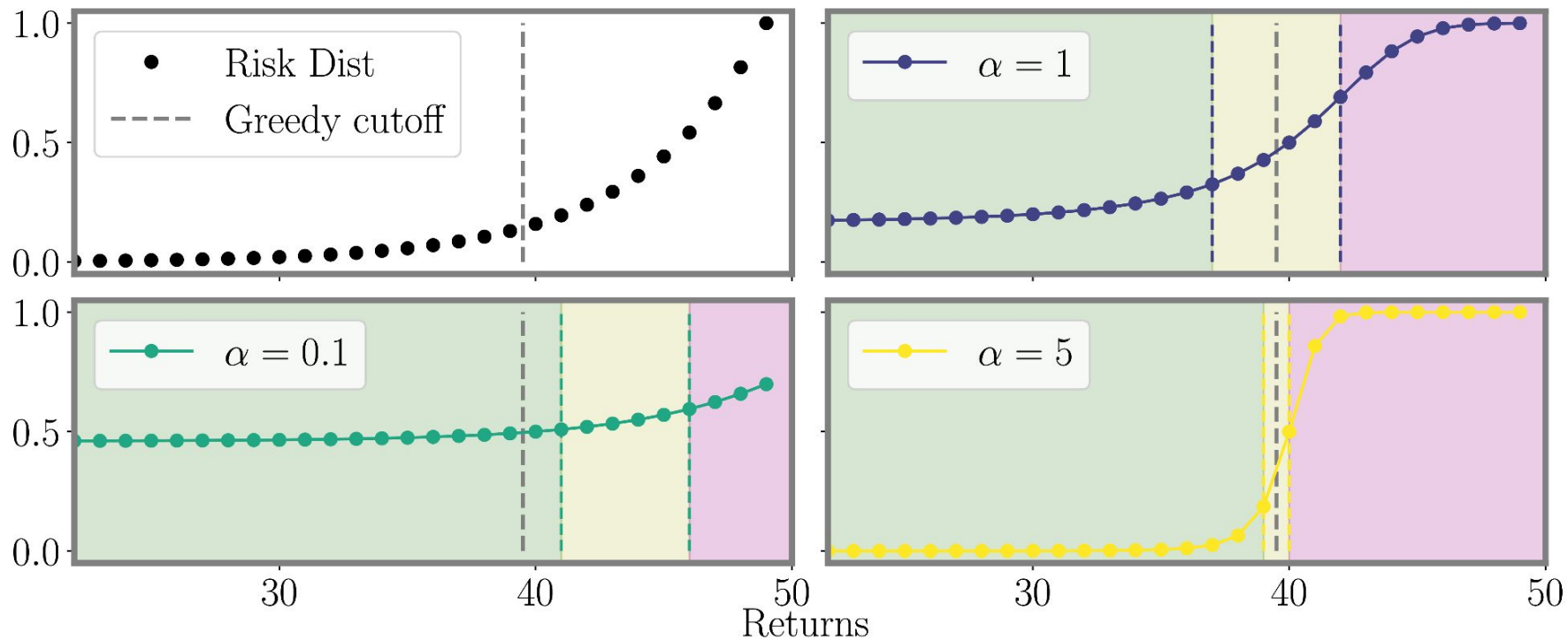
If only use this:

Unbiased estimate,  
but sub-optimal and low reward

# Adaptive Bin Sampling



# Adaptive Bin Sampling



# Adaptive Bin Sampling

Horvitz-Thompson estimator gives **unbiased** estimate.

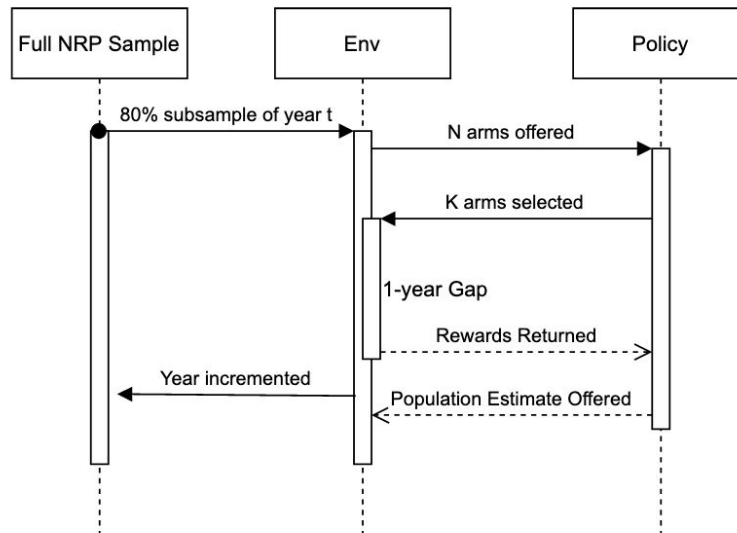
And we have **fine-grained** control over reward-variance trade-off.

$$\hat{\mu}_{HT}(t) = \frac{1}{\sum_a w_a} \sum_{a \in \mathcal{K}} \frac{w_a r_a}{p_a},$$

# Experiments

For NRP data years 2006-2014

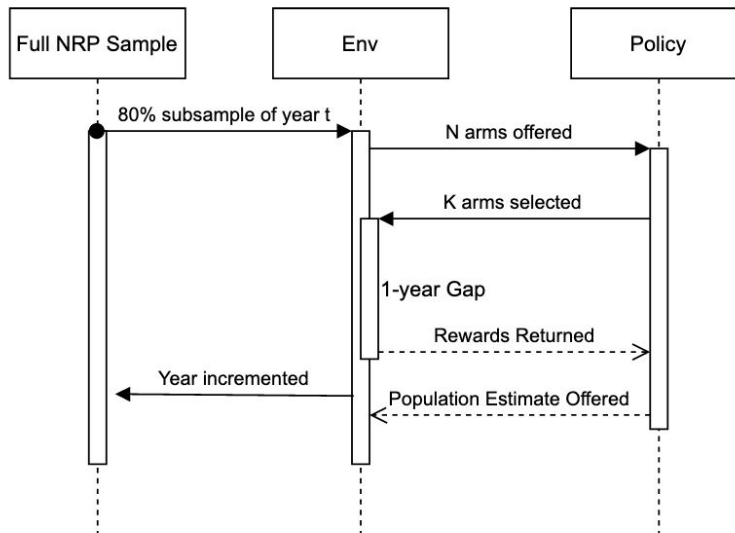
1. Take 80% subsample



# Experiments

For NRP data years 2006-2014

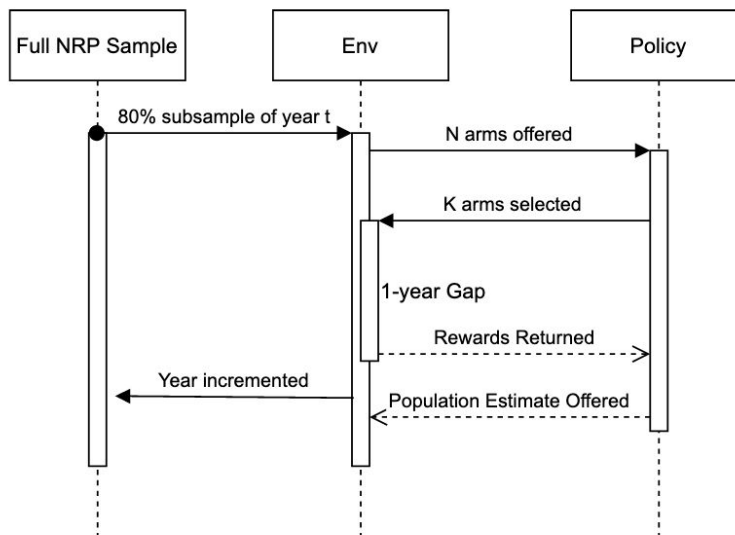
1. Take 80% subsample
2. Give selection policy  $\sim 500$  covariates from tax return data for each “arm” (tax return) in the sample



# Experiments

For NRP data years 2006-2014

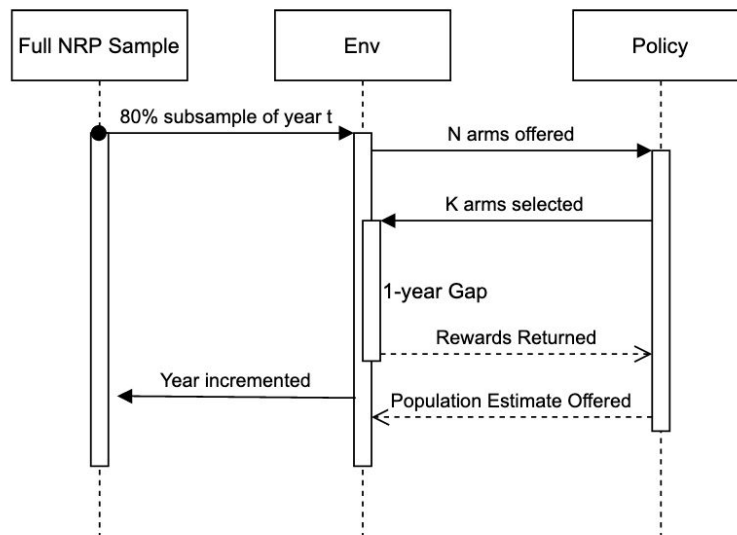
1. Take 80% subsample
2. Give selection policy ~500 covariates from tax return data for each “arm” (tax return) in the sample
3. Selection policy returns arms to audit



# Experiments

For NRP data years 2006-2014

1. Take 80% subsample
2. Give selection policy ~500 covariates from tax return data for each “arm” (tax return) in the sample
3. Selection policy returns arms to audit
4. Simulate a 1 year gap

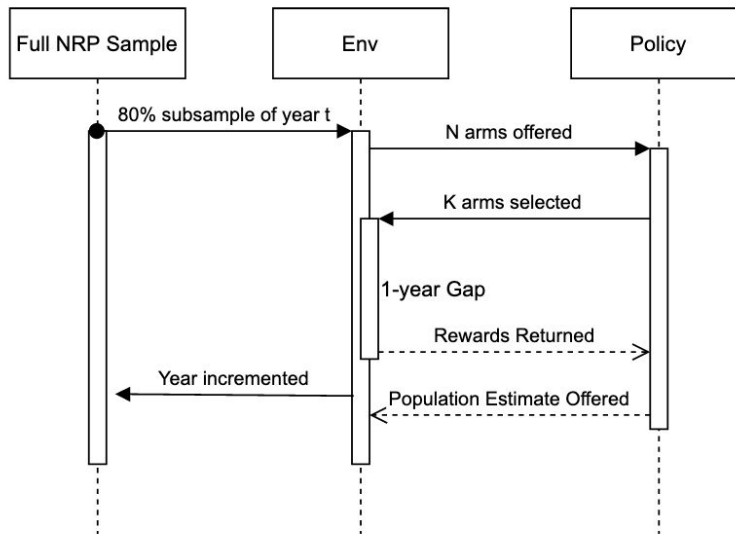




# Experiments

For NRP data years 2006-2014

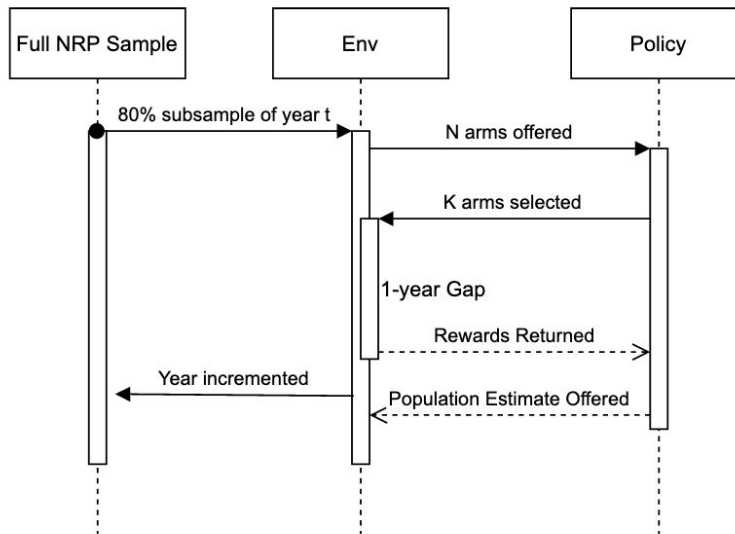
1. Take 80% subsample
2. Give selection policy  $\sim 500$  covariates from tax return data for each “arm” (tax return) in the sample
3. Selection policy returns arms to audit
4. Simulate a 1 year gap
5. Return the tax adjustment (reward) after that gap



# Experiments

For NRP data years 2006-2014

1. Take 80% subsample
2. Give selection policy  $\sim 500$  covariates from tax return data for each “arm” (tax return) in the sample
3. Selection policy returns arms to audit
4. Simulate a 1 year gap
5. Return the tax adjustment (reward) after that gap
6. Policy makes population estimate




# Experiments

<b>Best Reward Settings</b>					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Unbiased Methods	ABS-1	<b>\$41.5M*</b>	<b>0.4</b> ✓	31.0	<b>37.6%</b>
	$\epsilon$ -only	\$41.3M*	4.3 ✓	37.4	38.3%
	ABS-2	\$40.5M*	0.6 ✓	24.5	38.3%
	Random	\$12.7M	1.5 ✓	<b>14.7</b>	53.1%

# Experiments

10% ( $\epsilon$ ) random sample,  
rest greedy



<b>Best Reward Settings</b>					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Unbiased Methods	ABS-1	<b>\$41.5M*</b>	<b>0.4</b> ✓	31.0	<b>37.6%</b>
	$\epsilon$ -only	\$41.3M*	4.3 ✓	37.4	38.3%
	ABS-2	\$40.5M*	0.6 ✓	24.5	38.3%
	Random	\$12.7M	1.5 ✓	<b>14.7</b>	53.1%

# Experiments

Best Reward Settings					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Unbiased Methods	ABS-1	<b>\$41.5M*</b>	<b>0.4</b> ✓	31.0	<b>37.6%</b>
	$\epsilon$ -only	\$41.3M*	4.3 ✓	37.4	38.3%
	ABS-2	\$40.5M*	0.6 ✓	24.5	38.3%
	Random	\$12.7M	1.5 ✓	<b>14.7</b>	53.1%

Fully random sample every year,  
rest greedy

# Experiments

ABS can yield lower variance, similar reward, lower no-change rate, and retain unbiasedness

Best Reward Settings					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Unbiased Methods	ABS-1	\$41.5M*	0.4 ✓	31.0	37.6%
	$\epsilon$ -only	\$41.3M*	4.3 ✓	37.4	38.3%
	ABS-2	\$40.5M*	0.6 ✓	24.5	38.3%
	Random	\$12.7M	1.5 ✓	14.7	53.1%

# Experiments

Greedy tends to perform well in highly stochastic low-sample regime (which matches our experimental setup). (Bastani et al., 2022 proved this recently.)

Best Reward Settings					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Biased Methods	Greedy	\$43.6M*	16.4 <b>X</b>	8.8	<b>36.5%</b>
	UCB-1	\$42.4M*	15.3 <b>X</b>	9.4	38.6%
	$\epsilon$ -Greedy	\$41.3M*	<b>6.1 X</b>	<b>7.5</b>	38.3%
	UCB-2	\$40.7M*	15.6 <b>X</b>	10.21	40.7%

# Experiments

Best Reward Settings					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Biased Methods	Greedy	\$43.6M*	16.4 <b>X</b>	8.8	36.5%
	UCB-1	\$42.4M*	15.3 <b>X</b>	9.4	38.6%
	$\epsilon$ -Greedy	\$41.3M*	6.1 <b>X</b>	7.5	38.3%
	UCB-2	\$40.7M*	15.6 <b>X</b>	10.21	40.7%

Use regression model for both selection and population estimate. Means biased prediction, but slightly more reward and lower variance



# Experiments

## Best Reward Settings

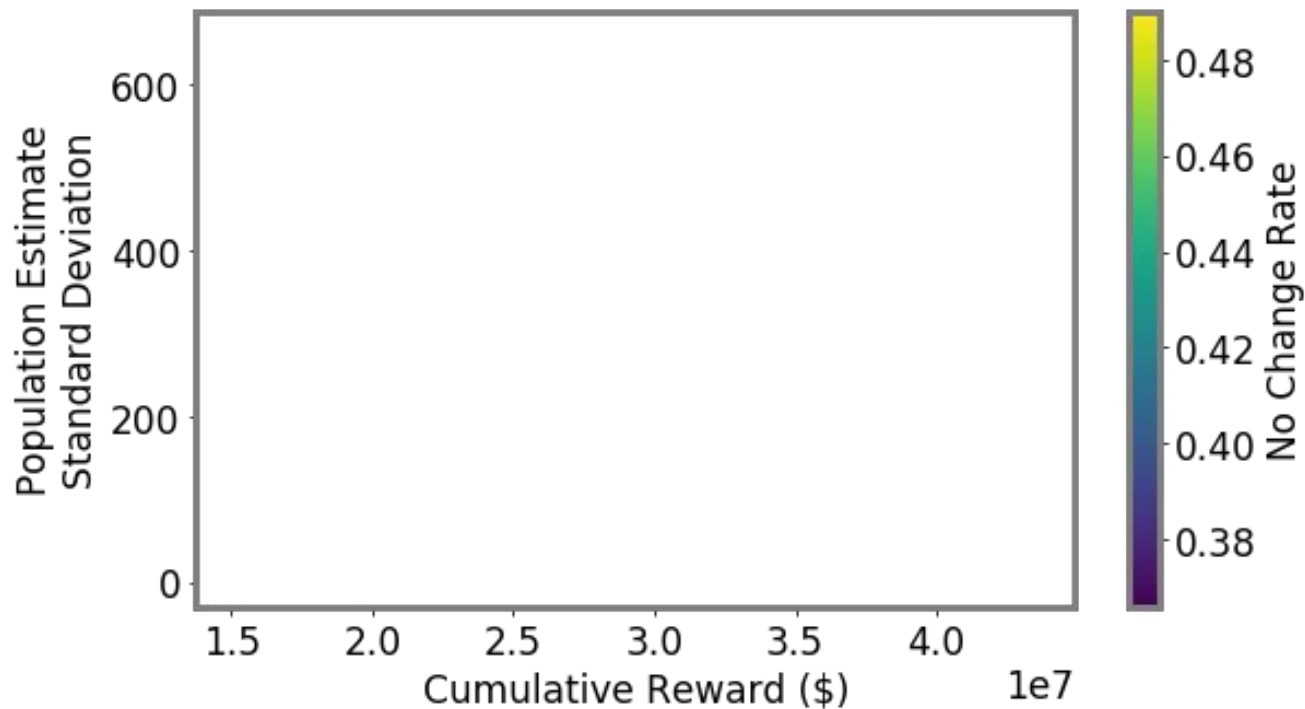
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Biased Methods	Greedy	\$43.6M*	16.4 <b>X</b>	8.8	<b>36.5%</b>
	UCB-1	\$42.4M*	15.3 <b>X</b>	9.4	38.6%
	$\epsilon$ -Greedy	\$41.3M*	<b>6.1 X</b>	<b>7.5</b>	38.3%
	UCB-2	\$40.7M*	15.6 <b>X</b>	10.21	40.7%

Even some randomness,  
reduces bias of model-based estimate,  
but not guaranteed.

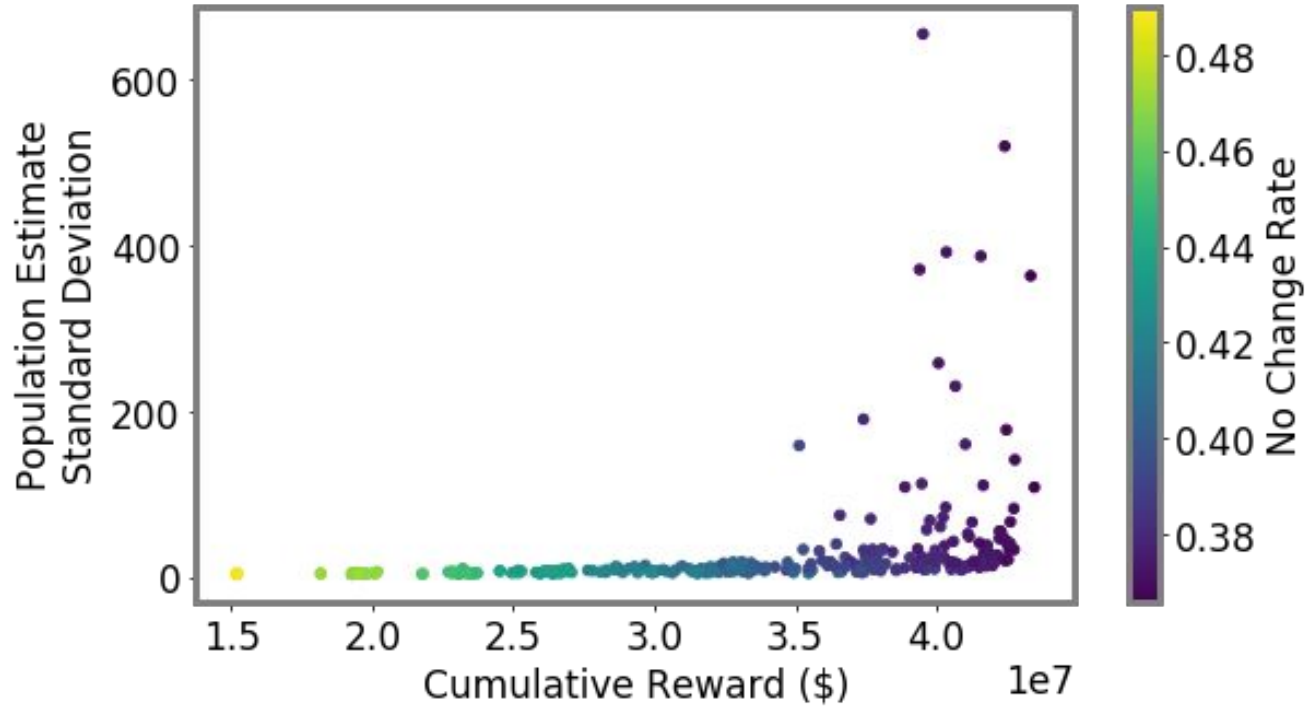
# Experiments

Best Reward Settings					
	<i>Policy</i>	<i>R</i>	$\mu_{PE}$	$\sigma_{PE}$	$\mu_{NR}$
Unbiased Methods	ABS-1	<b>\$41.5M*</b>	<b>0.4</b> ✓	31.0	<b>37.6%</b>
	$\epsilon$ -only	\$41.3M*	4.3 ✓	37.4	38.3%
	ABS-2	\$40.5M*	0.6 ✓	24.5	38.3%
	Random	\$12.7M	1.5 ✓	<b>14.7</b>	53.1%
Biased Methods	Greedy	<b>\$43.6M*</b>	16.4 ✗	8.8	<b>36.5%</b>
	UCB-1	\$42.4M*	15.3 ✗	9.4	38.6%
	$\epsilon$ -Greedy	\$41.3M*	<b>6.1</b> ✗	<b>7.5</b>	38.3%
	UCB-2	\$40.7M*	15.6 ✗	10.21	40.7%

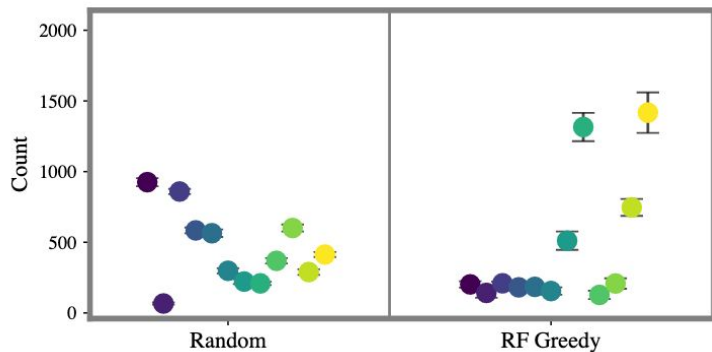
# ABS Enables Formal Tradeoff Between Precision and Reward



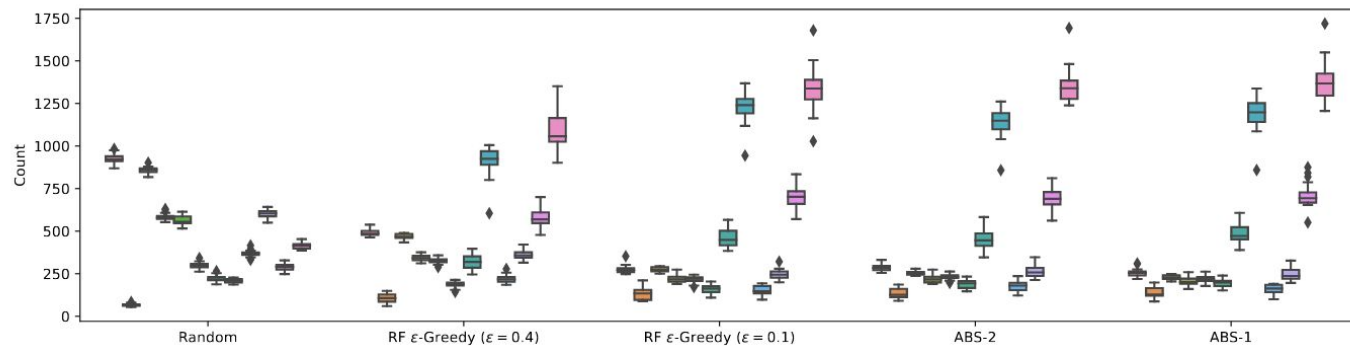
# ABS Enables Formal Tradeoff Between Precision and Reward



# More optimal methods sample higher incomes

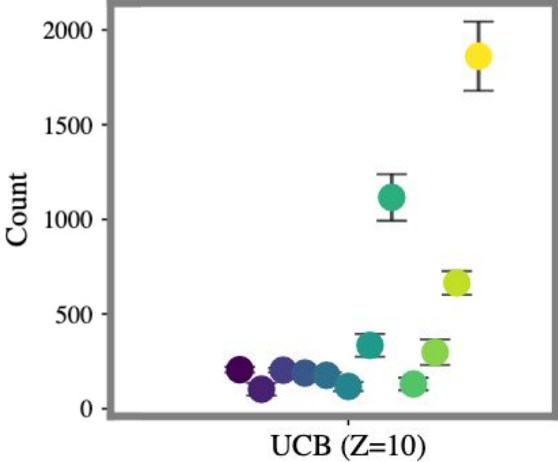


	EITC	TPI	C/F	E/2016	Farm
●	Y	<\$200K	<\$25K	N	N
●	Y	<\$200K	>\$25K	N	N
●	N	<\$200K	N	N	N
●	N	<\$200K	N	Y	N
●	N	<\$200K	<\$25K	N	N
●	N	<\$200K	\$25K-100K	N	N
●	N	<\$200K	\$100K-200K	N	N
●	N	<\$200K	>\$200K	N	N
●	N	<\$200K	N	N	Y
●	N	\$200K-\$1M	N	N	N
●	N	\$200K-\$1M	Y	N	N
●	N	>\$1M	N	N	N



■	EITC, TPI < \$200k, Sch C/F < \$25K
■	EITC, TPI < \$200k, Sch C/F > \$25K
■	No EITC, TPI < \$200k, no Sch C/E/F or F2106
■	No EITC, TPI < \$200k, Sch E or F2106, no Sch C/F
■	No EITC, TPI < \$200k, No Farm, Sch C/F < \$25k
■	No EITC, TPI < \$200k, No Farm, Sch C/F \$25k-100k
■	No EITC, TPI < \$200k, No Farm, Sch C/F \$100k-200k
■	No EITC, TPI < \$200k, No Farm, Sch C/F > \$200k
■	No EITC, TPI < \$200k, Farm
■	No EITC, TPI \$200k-\$1M, No Sch C/F
■	No EITC, TPI \$200k-\$1M, Sch C/F
■	No EITC, TPI > \$1M

# But heteroskedasticity can also drive sampling higher incomes



	EITC	TPI	C/F	E/2016	Farm
●	Y	<\$200K	<\$25K	N	N
●	Y	<\$200K	>\$25K	N	N
●	N	<\$200K	N	N	N
●	N	<\$200K	N	Y	N
●	N	<\$200K	<\$25K	N	N
●	N	<\$200K	\$25K-100K	N	N
●	N	<\$200K	\$100K-200K	N	N
●	N	<\$200K	>\$200K	N	N
●	N	<\$200K	N	N	Y
●	N	\$200K-\$1M	N	N	N
●	N	\$200K-\$1M	Y	N	N
●	N	>\$1M	N	N	N

# Takeaways

1. Unbiased estimation of population (e.g., average misreporting) **can still yield returns almost as high as greedy selection, with careful sampling and HT estimation.**
  - a. Suggests that a **unified optimize-and-estimate program** could be better and be more efficiently optimized.
2. Model-based population mechanisms are not guaranteed to be unbiased, but bias in practice can be reduced with some randomness.
3. More optimal methods tend to sample higher incomes in our experiments.
4. But heteroskedasticity also drives sampling of higher-incomes in uncertainty-based methods.