

Avicaching: Reducing Bias in Citizen Science with a Scalable Principle-Agent Game

Introduction

Citizen science programs outsource large-scale data collection, thus enabling researchers to expedite scientific discovery and environmental conservation.



Figure 1. A flow diagram for *eBird*, a citizen science program by the Cornell Lab of Ornithology.

Bias in Citizen Science

Although citizens assist in large-scale data collection, they collect data according to their own motivations, rather than systematically providing scientific observations. These biases introduce sampling bias, and manifest in the form of undesirable spatio-temporal clustering of submitted observations.

This sampling bias can be tackled with Principal-agent games. In these games, researchers (principal) deploy incentives, which could motivate citizens (agents) to complete the most crucial tasks, by factoring in the citizens' reasoning process. For example, *eBird* recently introduced the *Avicaching* game to reward citizen scientists for visiting Figure 2. Spatial clustering in *eBird* submissions previously under-sampled locations.



before 2014 in the US.



Figure 3. A Principal-agent game. The principal incentivizes agents to complete crucial tasks.

Anmol Kabra*¹

¹ak2426@cornell.edu

²yexiang@purdue.edu

Self-interest

Reducing Spatial Clustering

The principal must factor in agents' preferences for tasks when allocating rewards.

In this two-stage problem, the problem of learning the agents' behavior, with respect to their preferences v and principal's rewards r, is embedded in the problem of distributing rewards to maximize the principal's utility U_p . The behavior and rewards are subject to constraints $B_a(v)$ and $B_p(r)$ respectively.

(Principal)	$r^* = \operatorname{argmax}_r$	$U_p(v^*$
	subject to	$B_p(r)$
	(Agents)	$v^* = a$

Deploy Rewards Identification Problem Learn Agents' Behavior Pricing Problem Redistribute Rewards r*

Problem Formulation

Identification Problem

Consider *eBird* as an example. We extract the number of agents' visits to a location *u* from the dataset before and after a reward $r_{t,u}$ was applied during time period *t*. The visit densities are denoted as $x_t, y_t \in \mathbb{R}^n$ respectively.

To learn how agents' visits changed during t, we can learn a mapping $P : x_t \mapsto y_t$ such that $Px_t \approx y_t$. The entry $p_{u,v}$ in *P* is the probability of agents' efforts shifting from location v to u; the entry is thus a function of environmental features f_u , reward $r_{t,u}$, and the distance of v and u, all transformed by g with parameters w.

$$p_{u,v} = \frac{\exp(g(f_u, d_{u,v}, r_{t,u}; \boldsymbol{w}))}{\sum_{u'} \exp(g(f_{u'}, d_{u',v}, r_{t,u'}; \boldsymbol{w}))} \qquad \boldsymbol{w}^* = \operatorname{argmin}_{\boldsymbol{w}} \sum_{t} \|\boldsymbol{y}_t - \boldsymbol{P}_{f,D,r_t;\boldsymbol{w}} \boldsymbol{x}_t\|_2^2$$
$$= \operatorname{softmax}_{v}(q)[u]$$

Pricing Problem

To reduce spatial clustering, the pricing problem finds rewards that minimize variance $\|\boldsymbol{y} - \overline{\boldsymbol{y}}\|_1$ in predicted visit densities y = Px, where P is evaluated at w^* and $r^* = \operatorname{argmin}_r \frac{1}{n} \|y - \overline{y}\|_1$ $x = \frac{x'}{\|x'\|_1}, x' = \sum_t x_t.$ Subject to r = P

Additionally, the rewards must be nonnegative, sum to at most budget \mathcal{R} , and can only be placed at Avicaching locations given by *a*, where $a_u = 1$ if *u* is an Avicaching location, and 0 otherwise.

Scaling Avicaching using Machine Learning

A previous study folded the identification problem as linear constraints into the pricing problem by defining g (in expression for $p_{u,v}$) as a linear transformation. The pricing problem was then solved using a Mixed-Integer Programming (MIP) solver, which scaled poorly as the number of locations in the game increased.

We generalize the formulation by using a deep neural network for *g*, and scale both problems efficiently by harnessing Graphical Processing Units (GPUs).

Yexiang Xue² Carla P. Gomes ³

³gomes@cs.cornell.edu

$$(v^*, r)$$

= $\operatorname{argmax}_{v} U_a(v, r)$
= $\operatorname{subject}$ to $B_a(v)$



subject to $\boldsymbol{y} = \boldsymbol{P}_{f,D,w^*;r}\boldsymbol{x}$ $\forall u, (1 - a_u)r_u = 0$ $\forall u, r_u \geq 0$ $\|\boldsymbol{r}\|_1 \leq \mathcal{R}$

Neural Network Architecture

We build the input as follows: F[v] is a $n \times m$ matrix such that F[v][u] = $[f_u, d_{u,v}, r_{t,u}]$, a vector of *m* features that could influence agents to shift from v to u. This input F[v] is then fed forward into the network, which repeatedly takes linear combinations and applies non-linear functions to obtain P[v], as shown below.



Figure 4. A 3-layer neural network for the identification problem. *s* is the softmax function, and a_i are non-linear "activation" functions. ηe_v is a regularizer ($e_v[u] = 1$ if v = u and 0 otherwise) that controls how much P[v] looks like e_v .

The neural network model is essentially a sequence of batch-matrix operations that can be parallelized on GPUs. We can learn w^* and r^* using gradient descent updates on the same backpropagation framework. As shown below, the learned *w*^{*} from the identification problem become inputs to the pricing problem, which tunes the rewards *r* to minimize the pricing problem objective.



We compare	e our GPU-run m
ous studies ((n = 116 location)

	Loss	Runtime (s)		2-layer	6-layer
Random	1.014		Learned	1.073	1.025
Random Forest	0.491	26.4	MIP	1.110	_
BFGS	0.374	507.3	Random	1.169	1.303
2-layer	0.366	48.0			_
6-layer	0.358	647.8	Table 2. Pricing prob	plem. Con $\ \boldsymbol{u} - \overline{\boldsymbol{u}} \ _1$	nparing th

Table 1. Identification problem. We dataset along t as 75-5-20 train-vali and report the test loss values for d models (lower the better). Compare state-of-the-art BFGS-based model, networks perform better while taking

In both problems, we perform better than the previous state of the art while taking less time; our models are more than 800x faster than MIP in the pricing problem.

References

[1] Anmol Kabra, Yexiang Xue, and Car GPU-accelerated Principal-Agent Ga Citizen Science. Submitted to ACM COMPASS 2019. The**Cornell**Lab of Ornithology

Figure 5. Tensor representations of the identification (left) and the pricing (right) models. Dark blue tensors are static and red are optimizable. On the top-left are directions of dimensions.

Results

nodels' performance and speed with those from previns, m = 35 total features, T = 182 time periods).

48.0	
647 8	Table 2. Pricing problem. Comparing the
017.0	normalized objective, $\frac{1}{n} \frac{\ \boldsymbol{y} - \overline{\boldsymbol{y}}\ _1}{\overline{\boldsymbol{y}}}$ (lower the better),
e split the	on different reward allocation schemes, we see
id-test sets	that 2- and 6-layer networks better optimize the
different	objective than the MIP formulation. The MIP
ed to the	formulation takes \geq 36,000 seconds to run, while
, our	the GPU-run 6-layer network completes 10,000
ing less time.	iterations in 44.45 seconds.

	Acknowledgements
rla P. Gomes. ame for Scalable	We are grateful for funding by National Science Foundation
	(Grant Number 0832782, 1522054, 1059284, 1356308), ARO
	grant W911-NF-14-1-0498, the Leon Levy Foundation and
	the Wolf Creek Foundation.