

# Data-Driven Agent-Based Modeling, with Application to Rooftop Solar Adoption

Zhang, H., Y. Vorobeychik, J. Letchford, and K. Lakkaraju (2016). Data-Driven Agent-Based Modeling, with Application to Rooftop Solar Adoption, *Journal of Autonomous Agents and Multiagent Systems*, 30:1023-1049.

Presented by Joseph Abolt, Michael Kovar

# Agenda

- Motivation for Data-Driven Agent-Based Modeling
- The Model Itself
- Comparison to State-of-the-Art; Performance
- Paper's Conclusions
- Presenter's Conclusions

# Motivation for DDABM

# Rooftop Solar Adoption

- Solar energy is a clean and renewable resource – it's good
- Incentive programs increase adoption rates
- Difficult to accurately model how well incentives perform
- “Big Data” presents new opportunities
- *Note: DDABM can be applied elsewhere in addition*

# Predictive Models Flawed

- Traditional, non-ABM models:
  - Fail to accurately capture stochasticity of data
  - Mimic the aggregate trend, but cannot predict effects of causal influences
  - Cannot accurately predict far into the future
- Current ABM models:
  - Not developed robustly
  - Validation often qualitative
  - Quantitative validation uses same data as for calibration

# DDABM Advances Current Techniques

- Individual behavior learned offline by machine learning
- Meaningfully quantifies uncertainty about predictions
- Substantially outperforms current models
- Provides quantitative assessment of performance

# Data-Driven Agent-Based Modeling

# Assumptions

## 1. Time is discrete

- Decisions at time  $t = N$  are only affected by decisions at time  $t < N$ , for all  $N$

## 2. Agents are homogenous

- Suppose  $h(x)$  is agent behavior contingent on state  $x$
- Agents use the same  $h$  function;  $x$  can include personal details
- Heterogeneity through state, not decision-making process

## 3. Individuals make independent decisions at each time $t$ conditional on state $x$



# Terminology

- Let  $i$  index agents and  $t$  index time (to some time horizon  $T$ )
- Let  $x_{i,t}$  represent state of  $i$ th agent at time  $t$
- Let  $y_{i,t}$  represent  $i$ th agent's response
  - For solar adoption, decision at time  $t$  to adopt (1) or not (0)
- Let all data  $D = \{(x_{i,t}, y_{i,t})\}_{i, t=0\dots T}$

# Step 1 – Data Preparation

- Select a time threshold  $T_c$
- Split  $D$  into past and present data, and future data
- Calibration data =  $D_c = \{(x_{i,t}, y_{i,t})\}_{i, t \leq T_c}$
- Validation data =  $D_v = \{(x_{i,t}, y_{i,t})\}_{i, t > T_c}$

## Step 2 – Learn Agent Behavior

- Learn model  $h$  on  $D_C$  such that  $y_{i,t} = h(x_{i,t})$
- Use cross-validation on  $D_C$  for model (e.g., feature) selection

## Step 3 – Prepare Model

- Instantiate agents using  $h$  in the ABM
- Initialize the ABM to state  $x_{i,T_C}$  for all artificial agents  $i$

## Step 4 – Validation

- Validate ABM against  $D_v$  by running the model forward from  $x_{T_C}$

# DDABM as Applied to Solar Adoption

# Data From CSI (California Solar Initiative)

- Includes:
  - Individual-level adoption characteristics of residential solar projects in San Diego county
  - Property assessment for entire San Diego county
  - Electricity utilization data for most of the San Diego county 12 months prior to system installation
  - System size, reported cost, incentive amount, purchased/leased, date of incentive reservation, date of installation
  - May 2007 – April 2013 (~6 years, ~8,500 adopters)

# Accounting for Own vs Lease

- Net Present Value calculations difficult for own vs lease
- Only care about adoption
- Probability of adoption  $p(x) = p_L(x_L) + p_O(x_O) - (p_L(x_L) * p_O(x_O))$



# Agent Behavior as Logistic Regression

- Designed Logistic Regression models for own versus lease
- Includes both economic impacts and peer effects
- $R^2$  values, p-values, or other measures of reliability not provided

# Ownership Logistic Regression Model

Predictor	Coefficient
(Intercept)	-10.45
Owner Occupied (binary)	1.23
Number of Installations within 1 Mile Radius	3.19e-03
Number of Installations within ¼ Mile Radius	7.05e-03
Lease Option Available (binary)	0.73
Winter (binary)	-0.61
Spring (binary)	-0.19
Summer (binary)	-0.37
Installation Density in Zipcode	82.02
NPV (Purchase)	9.74e-06

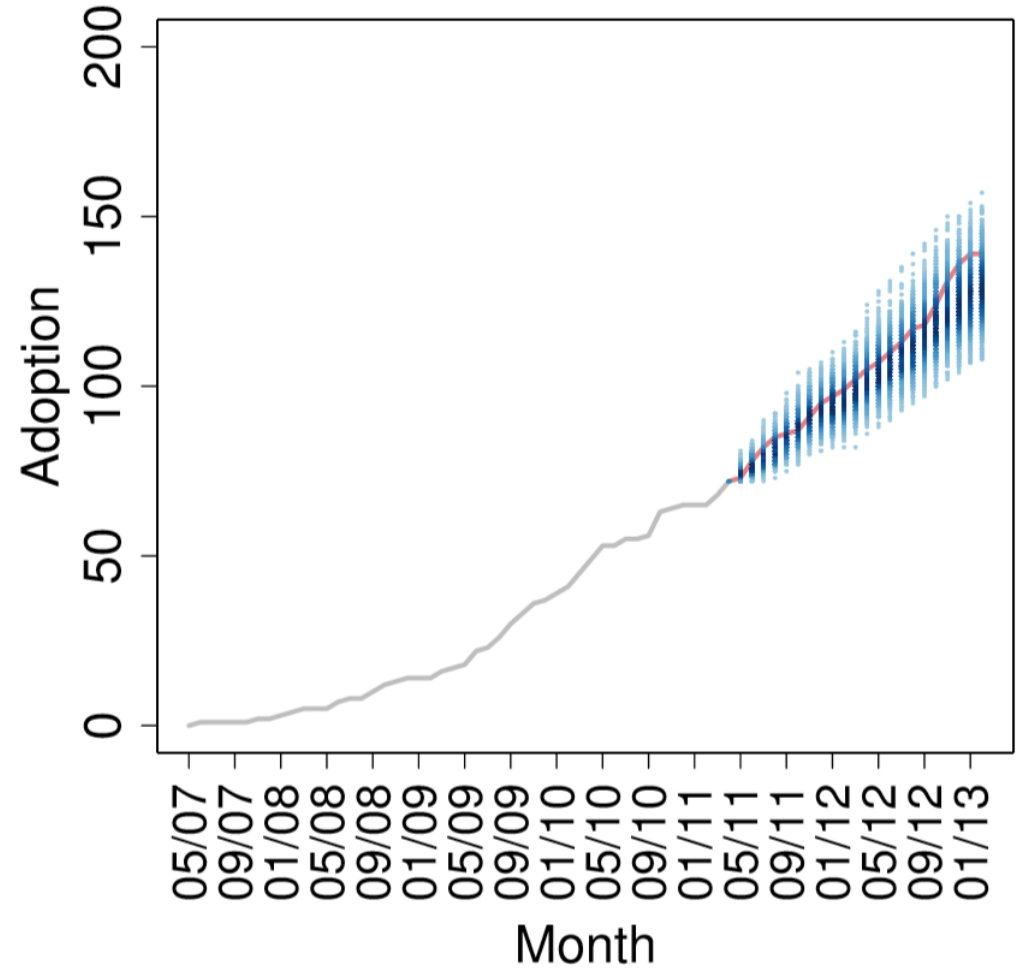
# Lease Logistic Regression Model

Predictor	Coefficient
(Intercept)	-14.04
Owner Occupied (binary)	1.00
Number of Installations within 2 Mile Radius	3.26e-03
Number of Installations within ¼ Mile Radius	9.58e-03
Lease Option Available (binary)	2.17
Winter (binary)	-0.40
Spring (binary)	0.30
Summer (binary)	-0.30
Installation Density in Zipcode	45.85
NPV (Lease)	1.03e-05

Performance,  
Comparison to State-of-the-Art

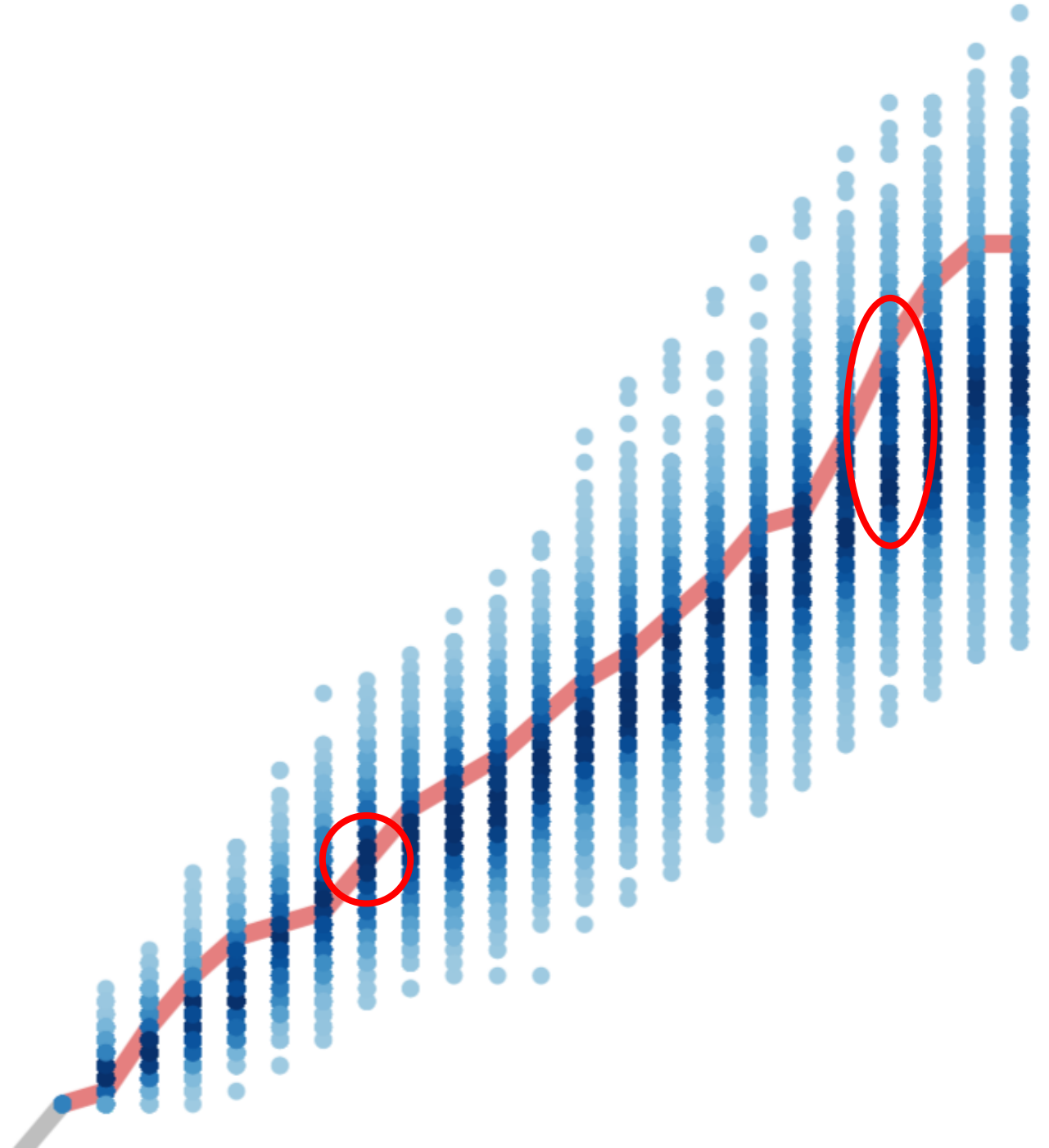
# Validation

- 1000 sample runs over a representative zip code (~13000 households)
- Training data through ~04/2011 (not explicitly stated)
- Testing data from ~04/2011 to 04/2013



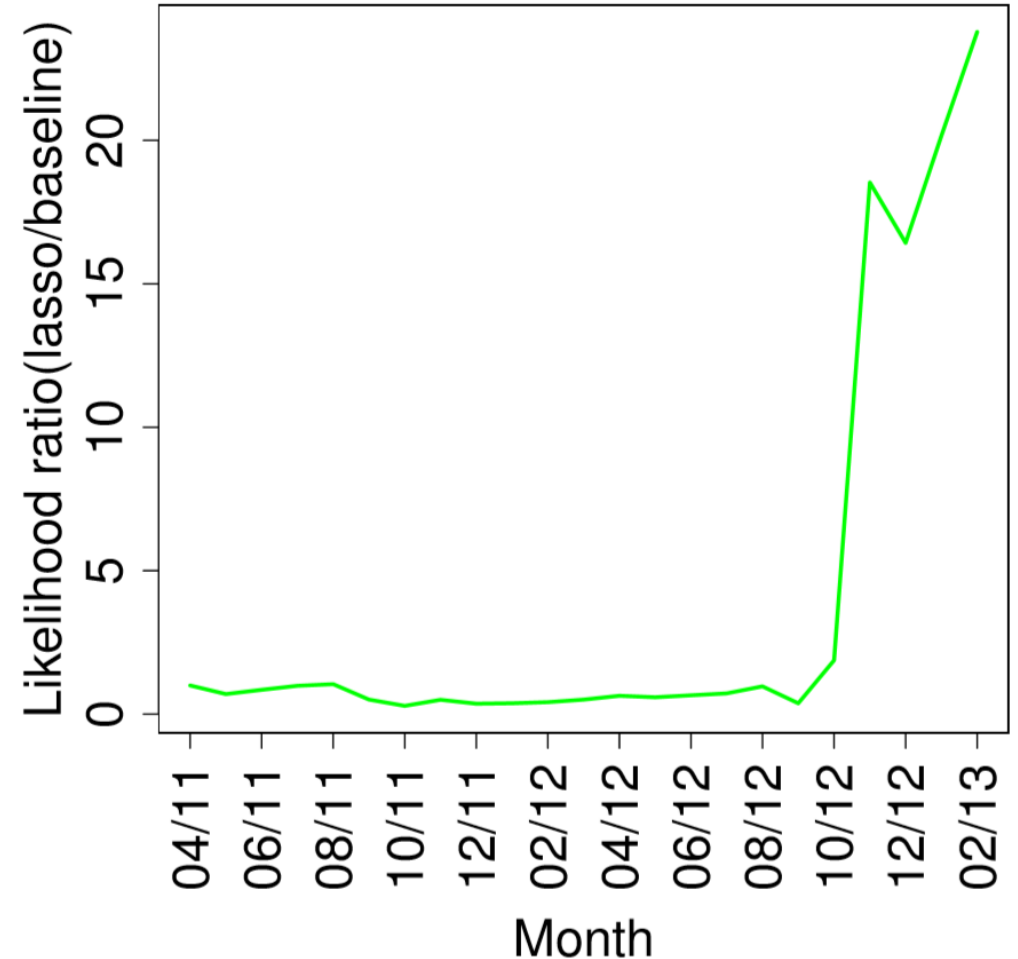
# Validation (cont.)

- True future passes through densest part of predictions
- Also note stochasticity yields normalized prediction ranges
- Can use stdev as measure of confidence



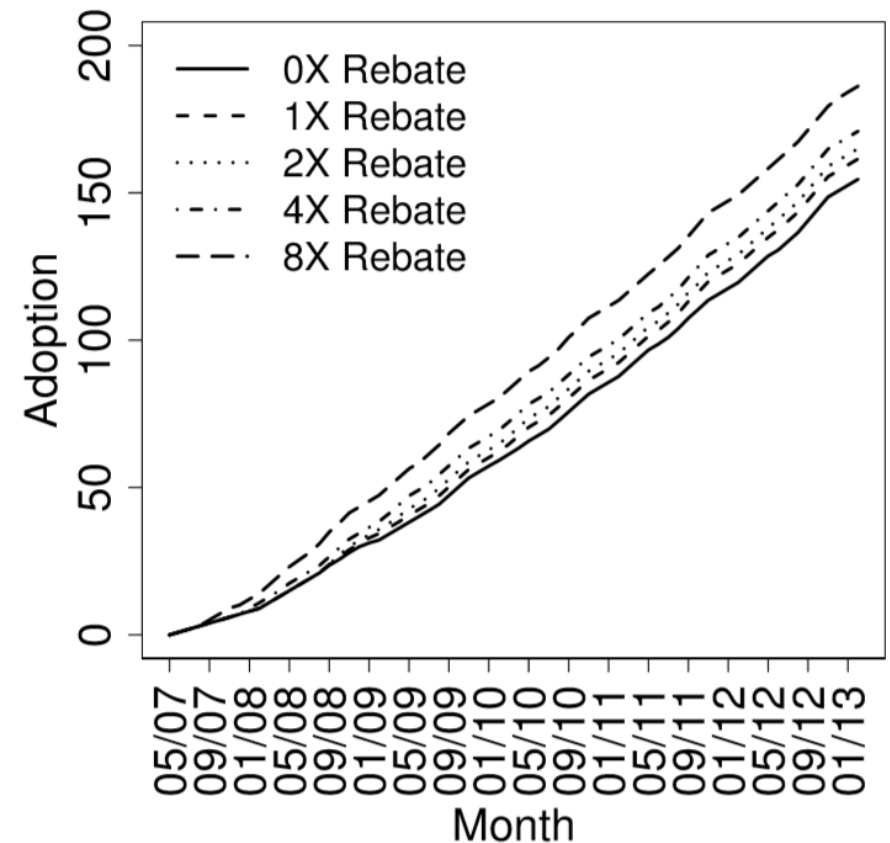
# Validation Against State-of-the-Art

- Similar to the state-of-the-art model in short-term
- At  $\sim 1.4$  years into the future, performance relative to state-of-the-art jumps an order of magnitude
- Cause: Predicting based on individuals versus aggregates



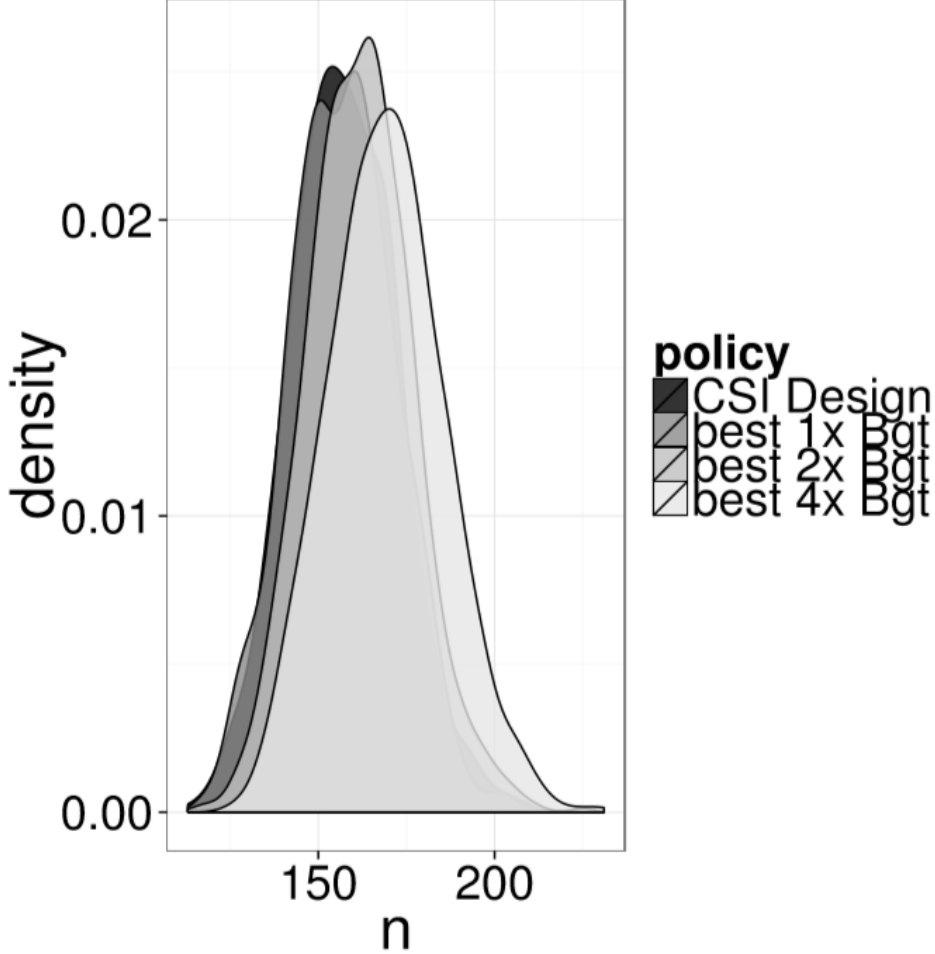
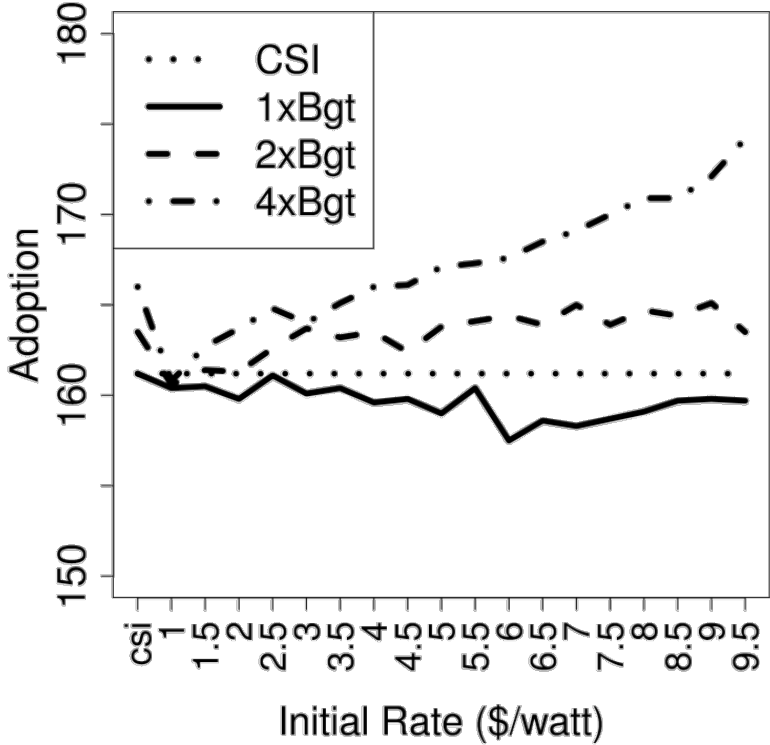
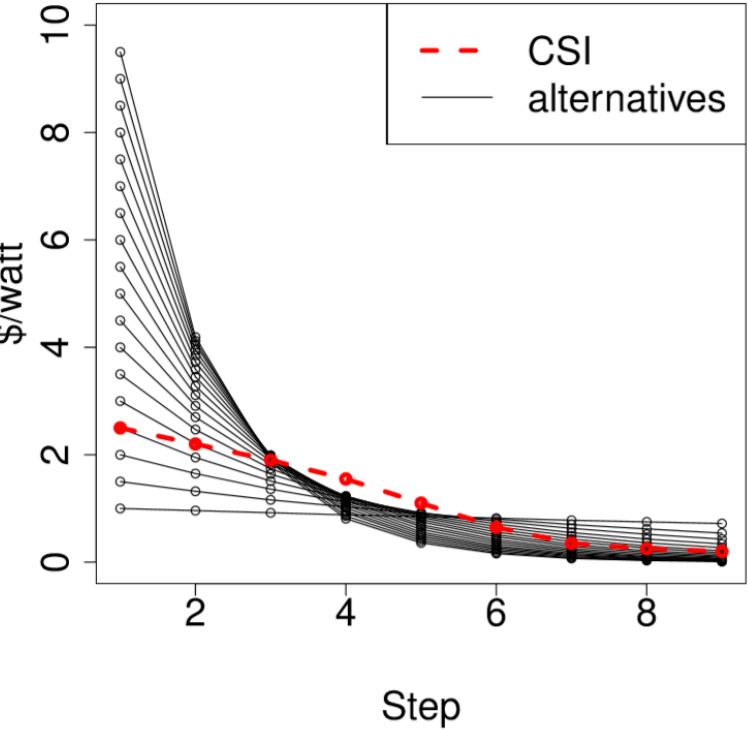
# Modeling Incentive Budgets

- Simplification instead of running 1000 models completely:
  - At each time  $t$ , generate  $t+1$  1000 times
  - Discard all but MLE
  - Reduces computational complexity
  - Maintains mean behavior
- Model versions of CSI incentive budget, from 0x (no rebate) to 8x
- Difference in adoption small



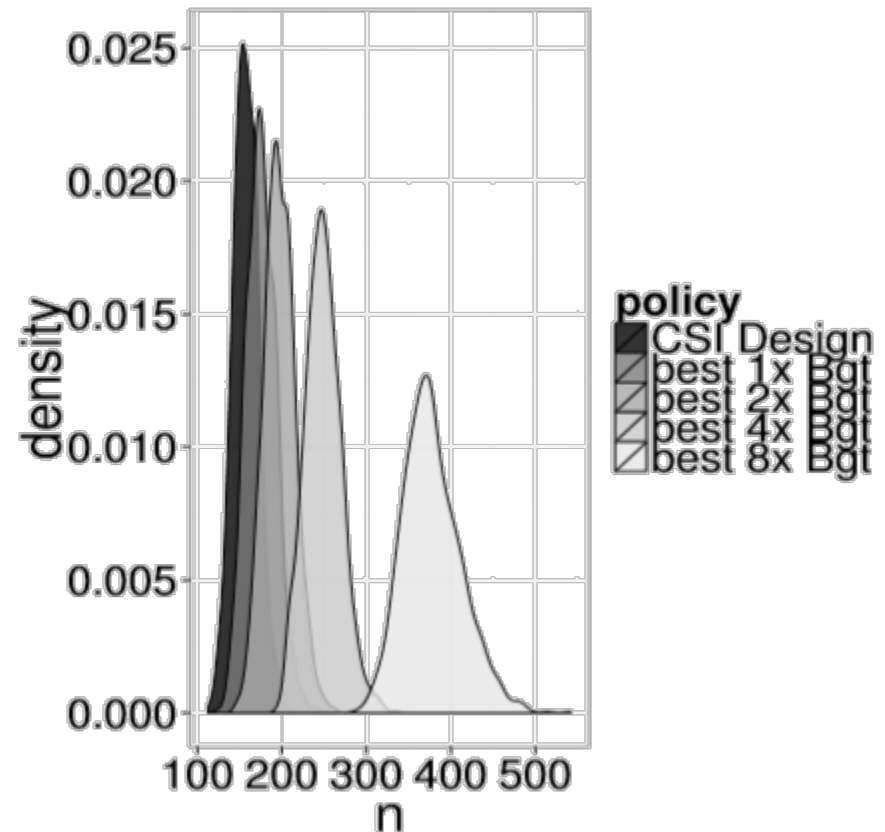


# Modeling Incentive Programs



# Modeling Seeding Programs

- Seeding programs take advantage of peer effects
- Incentive to seed early:
  - let peer effects last longer
- Incentive to seed later:
  - can produce more “seeds” for same budget



Conclusions from Paper

# Incentives Programs

- Significantly greater adoption possible from seeding programs
- Seeding more responsive to budget increase
  
- DDABM can provide quantitative estimates of adoption
- DDABM can provide confidence of estimates

# DDABM Viability

- Sufficient data features available to design DDABM
- Many applications beyond solar adoption
  
- Better than state-of-the-art by a magnitude
- Will improve with better data
- May improve with more sophisticated individual models

# DDABM Core Advancements

- Quantifiable verification of performance
- Quantifiable confidence measures
- Verification data temporally beyond calibration data
  - As opposed to reusing calibration data for validation

Conclusions from Presenters

# State-of-the-Art Comparison Flawed

- State-of-the-Art model built on different dataset
  - Used more data, because more data available in Italy
- Adapting to San Diego led to a double assumption
- Result: conflating income utility with square feet of house
- Additionally: data drawn from historical mean home sale prices
  - Not individual to agents, nor necessarily representative
- Also assumed proximity valid predictor of socioeconomic status
- Calls into doubt observed improvements over State-of-the-Art



# Reliability of $h$ in Question

- $h$  is the model for agent behavior
- Essentially a logistic regression model
- No analysis of the accuracy of  $h$  compared to calibration data provided
- No analysis of the features used in  $h$  provided
- Expected backwards-elimination of parameters and an  $R^2$  value

# Model's Confidence Unstated

- As a stochastic Bernoulli process, output at time  $t$  should be a Normal distribution
- No analysis of how broad 1 standard deviation is
  - From graphs, by the time DDABM surpasses SotA, mean  $\sim 125$ , range  $\sim 50$
- Unlikely all 1000 results were graphed, uncertain how representative the graph is
- If stdev = 25 around a mean of 125, confidence intervals prohibitively large

# Model's Usefulness Limited

- DDABM is only an improvement about 1.4 years out
- Therefore, only superior on problems 1.4 years out
  
- For shorter-term problems, SotA may be more computationally efficient
- Especially concerning when SotA crippled in implementation (mentioned earlier)

# Seeding-Incentive Comparison Misleading

- Seeding = giving people solar units
- Incentive = paying people to install solar units
  
- Seeding produces more adopters than incentives
- Uncertain if (adopters – seeded adopters) > adopters under incentive
  
- Similarly explains why seeding more responsive to budget; may be unrelated to model behavior

# DDABM Advantages Valid

- More robust evaluation of models to separate calibration and validation
  - Sets up ABM for bagging, and thence for boosting
  - Allows XGBoost-level behavior on complex problems
- Confidence measures vital for real-world application
- Confidence measures available; simply not demonstrated in paper

# Summary

- Much of the work lacks data to prove its validity
- The conceptual advancements are a solid foundation
  - Also pave the way for additional advancements
- Solar adoption proves good sample problem

Questions