# Modeling Nonresponse in Establishment Surveys: Using an Ensemble Tree Model to Create Nonresponse Propensity Scores and Detect Potential Bias in an Agricultural Survey

*Morgan Earp[1], Melissa Mitchell[2], Jaki McCarthy[3], and Frauke Kreuter[4]*

**BLS**

BUREAU OF LABOR STATISTICS
U.S. DEPARTMENT OF LABOR

www.bls.gov

# **Purpose**

- Developed classification trees to identify hardcore nonrespondents

- Assessed relationship between classification tree nonresponse propensity and actual nonresponse

- Created 10 classes based on classification tree nonresponse propensities to assess and compare nonresponse bias

# Motivation

- Attempt to reduce nonresponse bias, by identifying and targeting influential nonrespondents prior to survey administration

BLS

# ARMS Nonresponse Rates

Table 1. ARMS response rates 2000–2008

| Year | Sample size | Response rate (%) |
| --- | --- | --- |
| 2000 | 17,903 | 63 |
| 2001 | 13,313 | 64 |
| 2002 | 18,219 | 74 |
| 2003 | 33,861 | 63 |
| 2004 | 33,908 | 68 |
| 2005 | 34,937 | 71 |
| 2006 | 34,203 | 68 |
| 2007 | 31,924 | 70 |
| 2008 | 36,388 | 66 |

(p. 703)

BLS

# Methods

- Used an ensemble of classification trees to identify likely nonrespondents

- Used nonresponse propensity deciles to classify nonrespondents and assessed bias using the relative difference of the mean

BLS

# Classification Trees

- A "data mining" technique which segments a dataset using a series of simple rules to maximize dichotomies

- Creates subsets of records exhibiting a higher percentage of the "target"(respondent or nonrespondent)

BLS

# Splitting Criteria

■ Optimal Splitting Criteria

▶ Significance Testing
  - Uses the $p$ value as the stopping rule after applying a Bonferroni adjustment to mitigate bias toward variables w/ many values
    • Interval ($F$ test)
    • Nominal (Chi-Square)

▶ Variance Reduction
  - Measures the reduction in entropy, after adjusting for ordinal differences
    • Ordinal (Entropy)

# Classification Tree Proxy Data

- Imported Census of Agriculture (COA) response history for the ARMS III 2000-2008 Samples ($n =$ 254,632)

- Imported and matched 2002 COA data to be used as proxies of these operations characteristics
  - ▶ 78% match rate for 2002

# Types of Proxy Data

- Proxy data included 70 COA variables significantly related to ARMS nonresponse

  ▶ Operator Demographics
  ▶ Farm Type
  ▶ Size
  ▶ Commodities Raised
  ▶ Expenses
  ▶ Location

# Example Tree

**ARMS III Matched Sample (Training Data)**

37%
n = 79,616

**Sum of Poultry Inventory Data**

< 4
38%
n = 71,644

≥ 4
30%
n = 7,972

**Total Value of Products Sold + Government Payments**

< $110,005
27%
n = 30,904

≥ $110,005
45%
n = 40,740

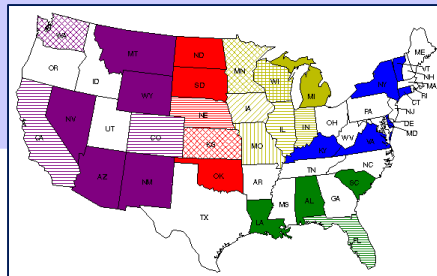**Total Sales – Not Under Production Contract (NUPC)**

< $844,879
41%
n = 31,211

≥ $844,879
57%
n = 9,529

**States**

AL, AZ, CA, CO, CT, DE, FL, IL, IN, IA, KS, KY, LA, MI, MN, MO, MT, NE, NV, NM, NY, ND, OK, SC, SD, VT, WA, & WY



Yes
70%
n = 3,346

No
52%
n = 1,118

BLS

# Analyses

- Assessed the relationship between classification propensity scores and nonresponse rates using logistic regression

- Assessed the relationship between classification propensity scores and nonresponse bias by plotting the relative bias of the mean by classification propensity score decile
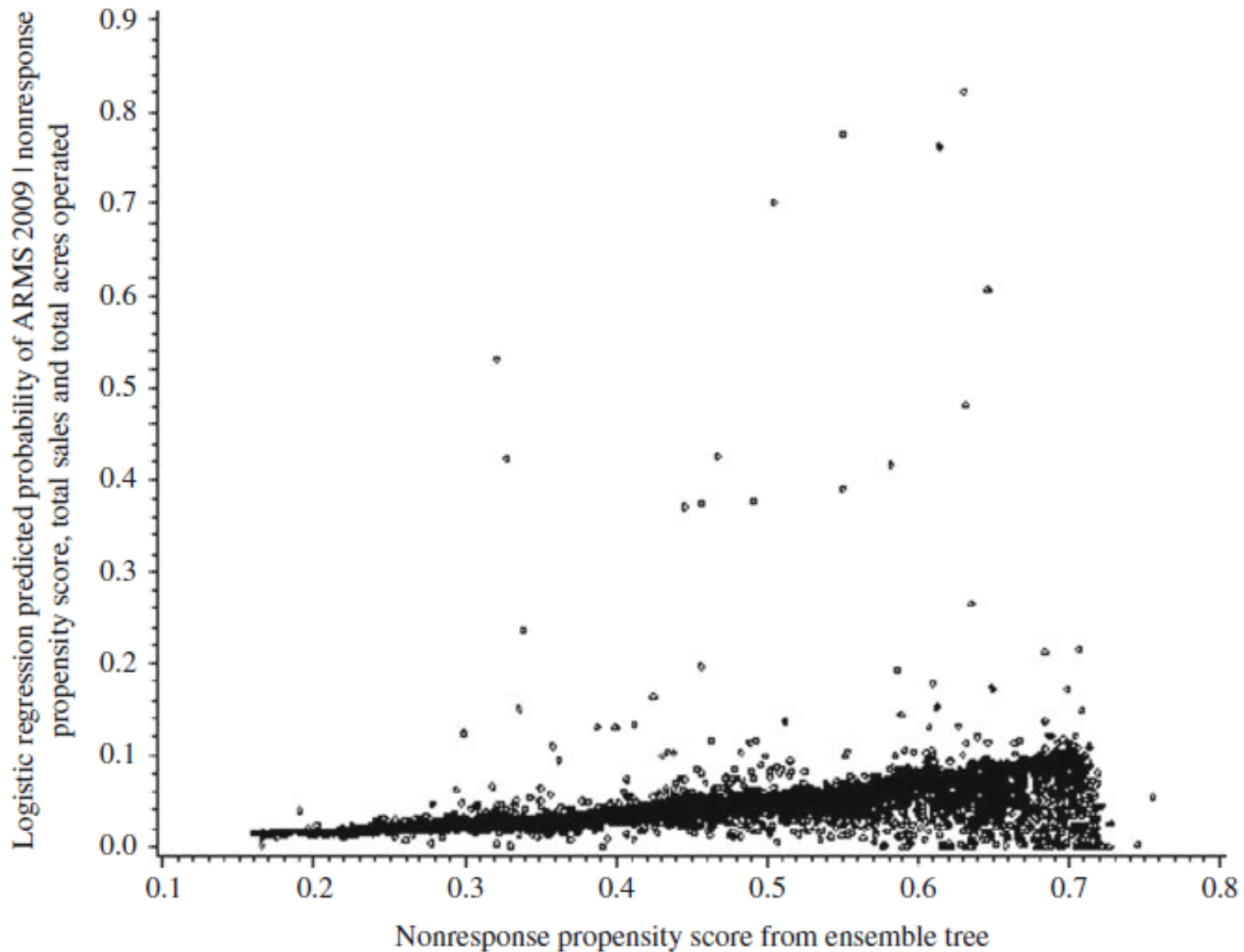
# **Variables**

- Inputs
  - Classification Tree Propensity Score
    - ARMS 2000-2008 nonresponse
    - Census 2002 operation characteristics

- Controls
  - Total Sales & Total Acres Operated
    - Census 2007

- Target
  - ARMS 2009 Nonresponse

BLS

Fig. 1. Plot of the logistic regression predicted probability of 2009 ARMS nonresponse given the ensemble tree nonresponse propensity score, 2007 total sales, and 2007 total acres operated, by the ensemble tree nonresponse propensity score

13

# Logistic Regression Results

Table 2.  *Logistic regression model fit statistics*

| Predictor | β | SE β | Wald's $\chi^2$ (df = 1) | p | $e^\beta$ Odds Ratio |
|---|---|---|---|---|---|
| | | Analysis of maximum likelihood estimates | | | |
| Constant | − 4.77 | 0.14 | 1191.55 | <.0001 | |
| Propensity score | 3.76 | .34 | 118.99 | <.0001 | 42.93 |
| Total sales | − 9.02−08 | 2.11E-08 | 18.35 | <.0001 | 1.00 |
| Total acres operated | 2.0E-05 | 3.19E-06 | 40.67 | <.0001 | 1.00 |

BLS

## Nonresponse rate by class

Nonresponse propensity calss from ensemble tree
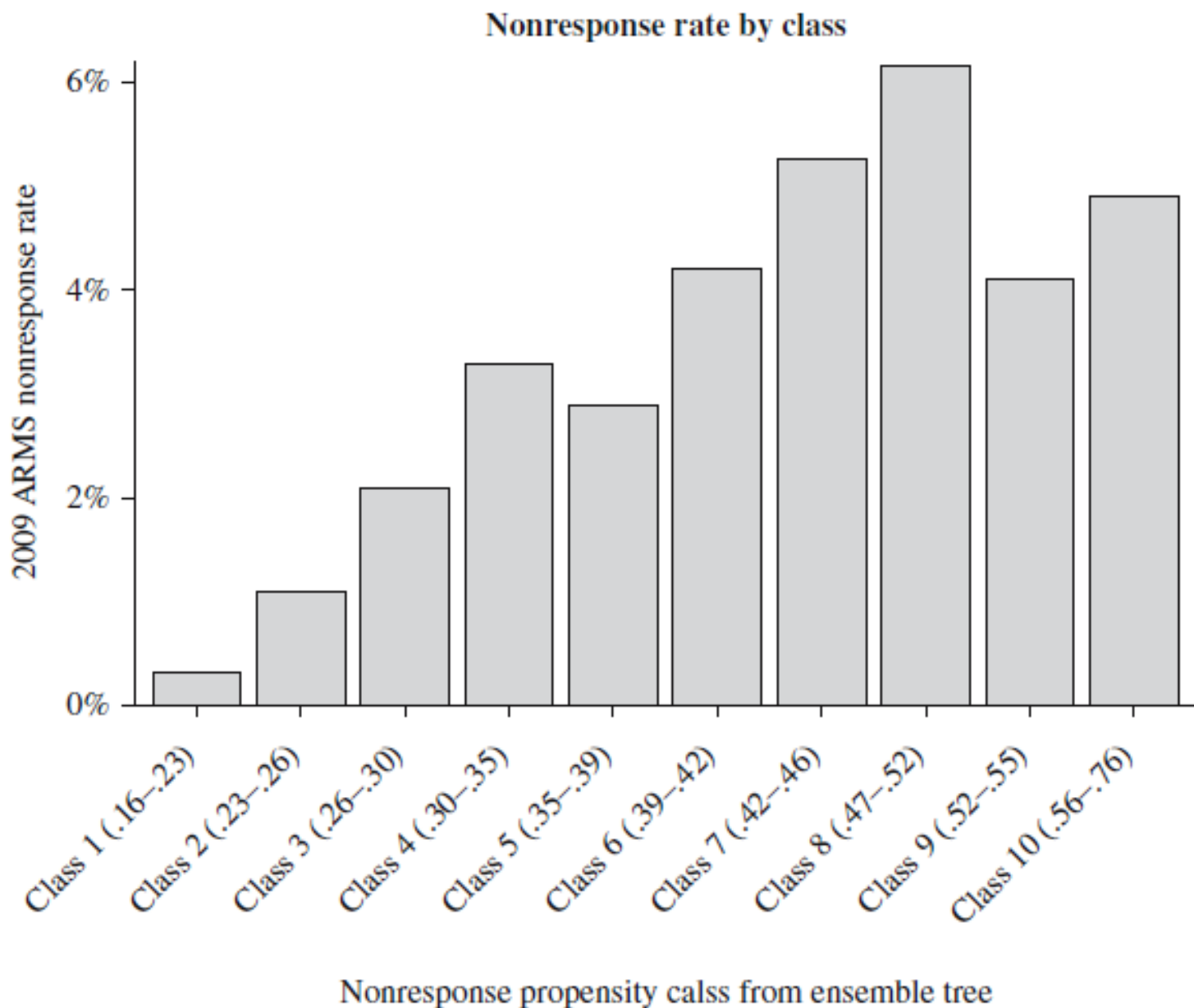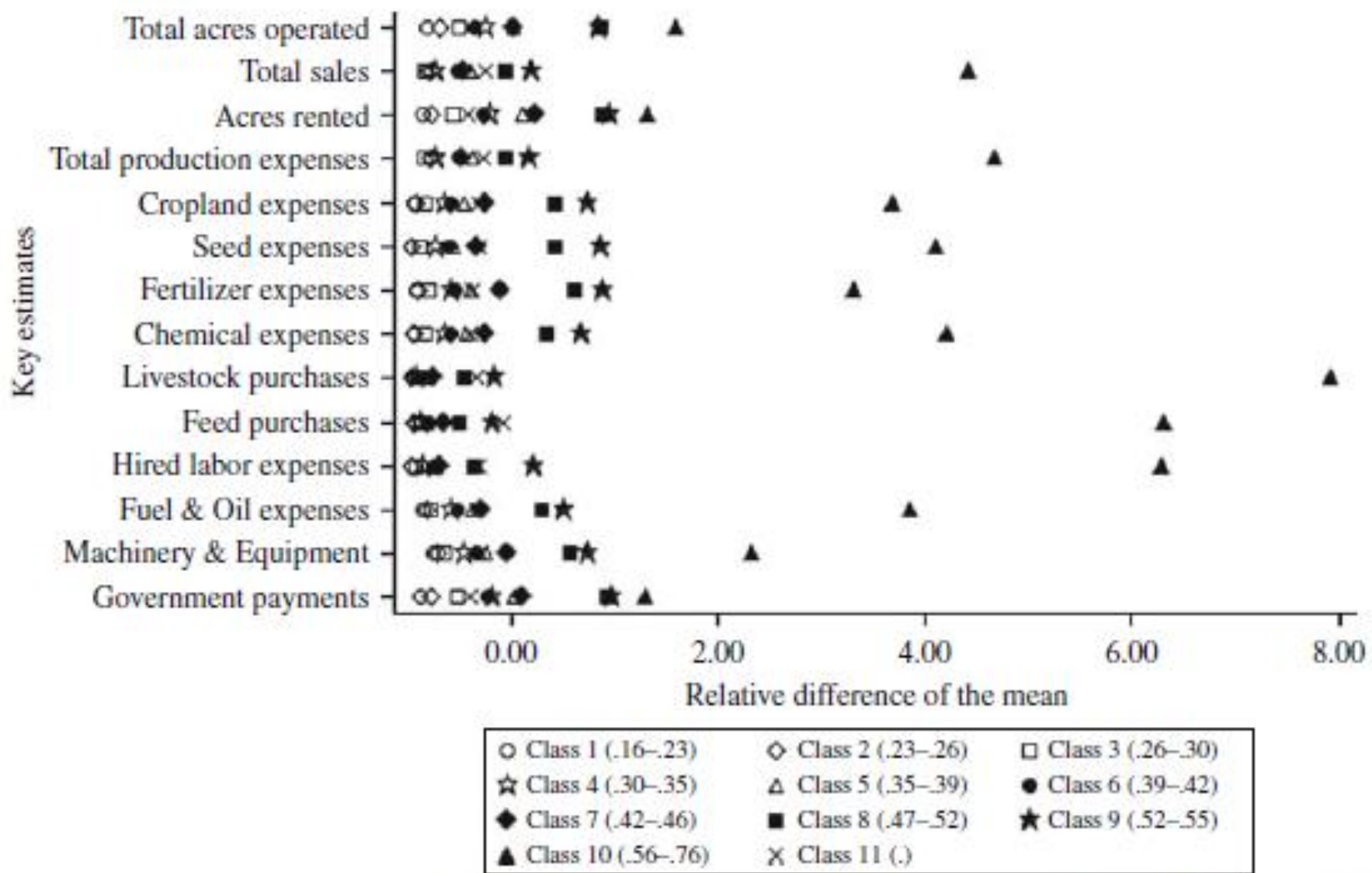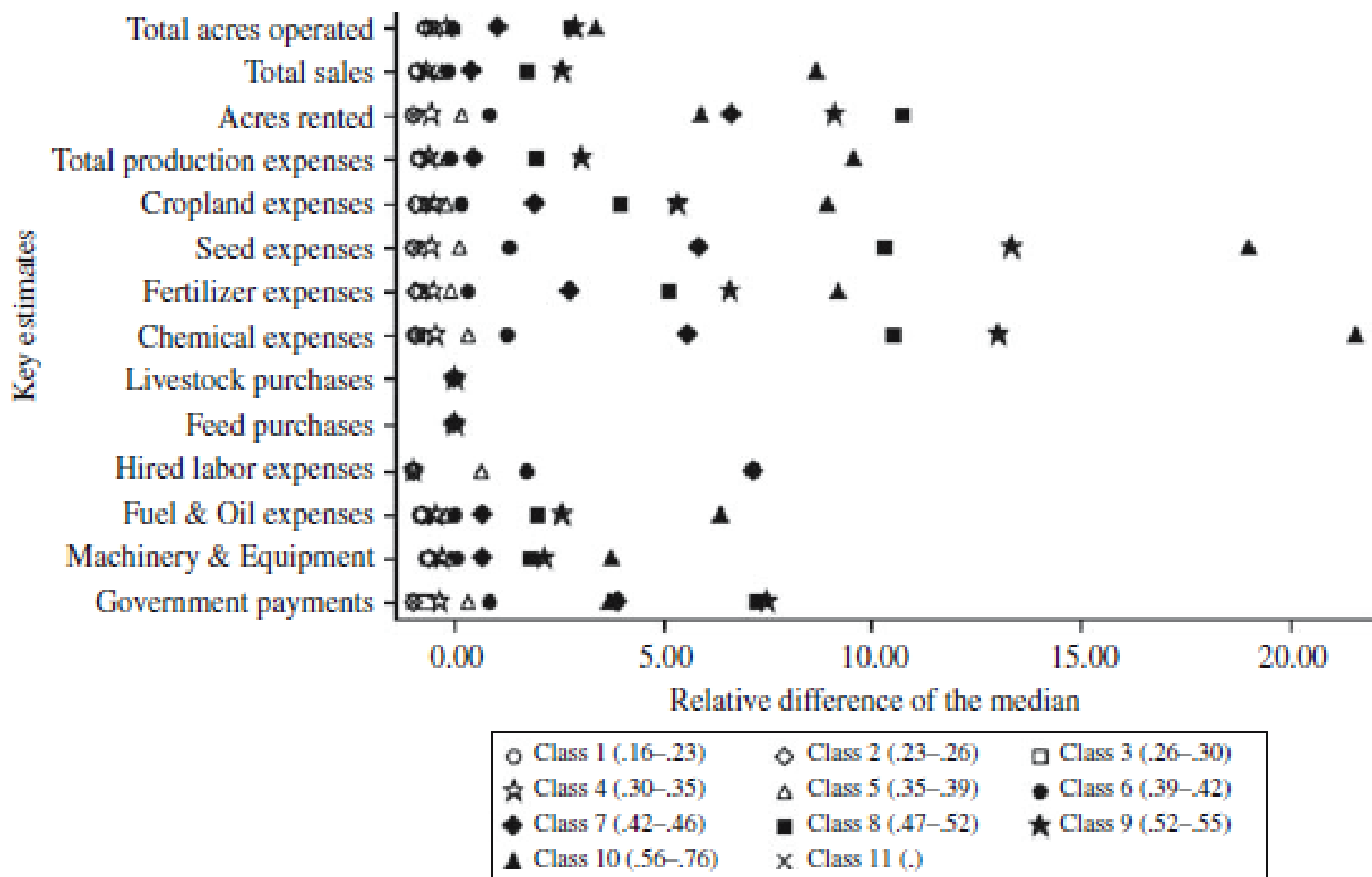
Fig. 2.    ARMS 2009 nonresponse rate by ensemble tree nonresponse propensity class

*Relative difference of the mean = [(class mean – overall mean)/overall mean]*

Fig. 3. Relative difference of the mean for key estimates by nonresponse propensity class

Fig. 4. Relative difference of the median for key estimates by nonresponse propensity class

Relative difference of the median = [(class median – overall median)/overall median]

# Conclusion

- Easily identify characteristics associate w/ nonresponse
- Can ensure that each variable is considered once in the overall average model
- These propensity scores were positively correlated with the amount of potential bias across several key estimates

BLS

# Conclusion

- We would like to compare this tree method w/ random forests

- They are currently being used to pre-score samples prior to data collection to ensure that those farms that are least likely to respond and most likely to bias estimates as a result receive special attention.

# Contact Information

Contact Information
Morgan Earp
Earp.Morgan@bls.gov

BLS
BUREAU OF LABOR STATISTICS
U.S. DEPARTMENT OF LABOR