



AgEcon SEARCH
RESEARCH IN AGRICULTURAL & APPLIED ECONOMICS

The World's Largest Open Access Agricultural & Applied Economics Digital Library

This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.

Help ensure our sustainability.

Give to AgEcon Search

AgEcon Search

<http://ageconsearch.umn.edu>

aesearch@umn.edu

*Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.*

No endorsement of AgEcon Search or its fundraising activities by the author(s) of the following work or their employer(s) is intended or implied.

Introduction

- Suppression of employment records in the US County Business Patterns (CBP) data sets constrains the detail of new methods and recent advances in the analysis of the geographic distribution of firms and employment.
- Data sets created by imputation procedures can be purchased, but cost often puts them beyond the reach of many research budgets.
- Fortunately, methods exist whereby researchers can impute suppressed employment records.
- A comparison of these procedures is necessary to assess the accuracy and flexibility of each.

Objectives

This research compares an existing goal programming approach created by Zhang and Guldmann (2009) (the ZG approach) and a new method, an Iteratively Constrained Rebalanced Matrix (ICRM) procedure with known, randomly suppressed data sets, constructed similar to the CBP.

Exploring the CBP

FIPS	NAICS	empflag	emp	est	n1_4	n5_9	n10_19
4001	22----	B	0	9	7	1	1

- Data is arranged in a hierarchy of geographical (FIPS) and industry (NAICS) classification.
- Suppressed record noted by letter in “empflag” column.
- Each employment flag corresponds with a range of employment in which the suppressed “emp” record falls.
- Establishment size ranges provide additional information about employment. In the example above, there are 9 establishments, 7 of which have between 1 and 4 employees, 1 has between 5 and 9 employees, and 1 has between 10 and 19. These intervals extend out to a range of 5000+ employees.
- From the flag and site interval information, the uncertainty about the suppressed employment value is narrowed from 20-99 employees, to 22-56 employees (57% decrease in uncertainty).

Methods

The ZG Approach (Zhang and Guldmann, 2009)

The goal programming approach proposed by ZG minimizes the sum of deviations of imputed estimates from constraints based on employment flag intervals, establishment size ranges, and known county and sector employment totals.

The ICRM Procedure

This approach uses nonlinear programming and a Gauss-Seidel minimization routine, in a process similar to the standard RAS algorithm. The procedure is completed iteratively, using the following steps:

1. The objective function minimizes the cross-entropy of row and column scaling factors,

$$\min_{y_i, z_j} \sum_i y_i \ln(y_i) + \sum_j z_j \ln(z_j)$$

subject to

$$\sum_j a_{0(ij)}^t \cdot y_i \cdot z_j = c_i \quad \forall j$$

$$\sum_i a_{0(ij)}^t \cdot y_i \cdot z_j = s_j \quad \forall i$$

2. The values for $a_{0(ij)}^t$ are imputed such that

$$a_{ij}^{t*} = \begin{cases} pmin_{ij} \leq a_{0(ij)}^{t-1} \cdot y_i^* \cdot z_j^* \leq pmax_{ij} & \text{if } a_{0(ij)} \text{ is suppressed} \\ a_{0(ij)} & \text{otherwise} \end{cases}$$

where

$$pmin_{ij} = MAX \left(fmin_{ij}, \sum_k tmin_k * t_{ijk} \right)$$

$$pmax_{ij} = MIN \left(fmax_{ij}, \sum_k tmax_k * t_{ijk} \right)$$

3. Convergence of the vectors is checked such that the difference in A^t and A^{t-1} after t iterations is a very small number.
4. If convergence is not reached, return to step one.

Identifier	Description
i	county index
j	sector index (1→21 for NAICS 11→99)
k	establishment size interval index (1→12)
t_{ijk}	number of establishments in county i , sector j , and establishment interval k
s_j	total employment for sector j across spatial units
c_i	total employment for county i across all sectors
$tmax_k$	employment for upper bound for interval k
$tmin_k$	employment for lower bound for interval k
$fmax_{ij}$	maximum employment of sector j in county i , if flagged
$fmin_{ij}$	minimum employment of sector j in county i , if flagged

Monte Carlo Simulation

- Simulated employment data sets that replicate the structure of a CBP data set.
- Increasing percentages of data were randomly suppressed in the simulated data sets.
- Goodness of fit was assessed for each method as the sum of the squared deviations between the true and imputed values to compare precision of the imputation procedures.

Results

- The ZG imputation procedure was feasible in 100% of the Monte Carlo trials.
- The unconstrained and constrained ICRM procedures had difficulty converging at lower levels of data suppression (<50%).
- All approaches were superior to the mid-point imputation approach.
- Compared to the “true” simulated employment patterns, mean absolute percent deviations of the constrained ICRM were lower than the ZG approach as the percent of data suppressed increased.

Conclusions and Further research

- Proposed method appears promising, computationally and in terms of precision.
- Method appears to work well when about half of the data are suppressed.
- Determine frequency of unconstrained ICRM estimated missing employment data outside relevant bounds.
- Incorporate ICRM bound constraints directly into nonlinear program.
- Investigate adaptive constraints, updating bounds at each iteration for ICRM approach.

Citations

- Zhang, S., and J.M. Guldmann. 2009. Estimating Suppressed Data in Regional Economic Databases: A Goal-Programming Approach. *European Journal of Operational Research* 192(2009): 521-537.