

WORKING P A P E R

Under-Reporting of Medicaid and Welfare in the Current Population Survey

JACOB ALEX KLERMAN, JEANNE S. RINGEL, AND
BETH ROTH

WR-169-3

March 2005

This product is part of the RAND Labor and Population working paper series. RAND working papers are intended to share researchers' latest findings and to solicit additional peer review. This paper has been peer reviewed but not edited. Unless otherwise indicated, working papers can be quoted and cited without permission of the author, provided the source is clearly referred to as a working paper. RAND's publications do not necessarily reflect the opinions of its research clients and sponsors. RAND® is a registered trademark.



LABOR AND POPULATION

Under-Reporting of Medicaid and Welfare in the Current Population Survey

Jacob Alex Klerman, Jeanne S. Ringel, and Beth Roth

Latest Revision: March, 2005

Address Correspondence to: Jacob Alex Klerman, RAND, 1700 Main St., Santa Monica, CA 90407, 1-310-393-0411, x289, Jacob_Klerman@rand.org.

This research was jointly funded by the California Health Care Foundation and the U.S. Department of Health and Human Services, Administration for Children and Families (Grant No. 01ASPE377A). The U.S. Bureau of the Census and the California Department of Health Services provided the data. The research in this paper was conducted while the authors were Research Associates at the Center for Economic Studies, U.S. Census Bureau. Research results and conclusions expressed are those of the authors and do not necessarily indicate concurrence by the U.S. Census Bureau, the Center for Economic Studies, the California Department of Health Services, RAND, or its sponsors including U.S. Department of Health and Human Services, Administration for Children and Families and the California Health Care Foundation.

This research has benefited greatly from the comments of Leonard Sternbach, and Richard Bavier at US DHHS; Ingrid Aguirre-Happoldt and Chris Perrone at the California Health Care Foundation; participants at a session at the Population Association of America Annual Meetings; and an anonymous RAND reviewer.

Preface

Conventional estimates of the number of uninsured Californians are derived from the Current Population Survey (CPS). Unfortunately, CPS estimates of the number of people receiving Medi-Cal and welfare (AFDC/CalWORKs) are well below the numbers implied by official Medi-Cal records, suggesting that the conventional estimates of the number of uninsured Californians (and their characteristics) are seriously flawed.

To improve our understanding of these issues, the California HealthCare Foundation (through its then separate the Medi-Cal Policy Institute—MCPI) and the U.S. Department of Health and Human Services, Administration for Children and Families (DHHS-ACF) funded RAND to match CPS data to individual-level administrative data for the Medi-Cal program. With the cooperation of the California Department of Health Services (CDHS), the U.S. Bureau of the Census, and the California Census Research Data Center (CCRDC), that match was performed. This document describes the findings of the analysis of those matched data.

Table of Contents

Preface	iii
Table of Contents	v
List of Tables	vii
List of Figures.....	ix
Summary.....	xi
The Current Population Survey (CPS), Under-Reporting, and Matching.....	xi
The Magnitude of Under-Reporting And Our Imputation Model.....	xii
The Effects of Under-Reporting on Estimates of Medi-Cal Enrollment Rates	xiv
The Effects of Under-Reporting on Estimates of Uninsurance	xiv
Summary.....	xv
Acknowledgments.....	xvii
Glossary, List of Symbols, etc.	xix
1. Introduction.....	1
Plan of the Report	1
Previous Literature on Under-Reporting	2
2. Medi-Cal, the CPS, The MEDS, and Under-Reporting.....	5
The Medi-Cal Program	5
The Medi-Cal Eligibility Data System (MEDS)	6
The Current Population Survey (CPS).....	8
Levels and Trends In Mis-Reporting	11
Discussion.....	12
3. The Matched Data.....	13
Matching	13
Simple Reporting Rates.....	17
Who Mis-reports?	19
“Behavioral Regressions” and “Imputational Regressions”	21
Understanding the Reporting Errors	22
CPS Reference Period.....	24
Time Trends.....	26
CPS Imputations	26
Conclusion	28
4. Extrapolating to the Full Data.....	29
The Identification Problem.....	29
Imputing the Data.....	32
Stratifying and Adjusting for Covariates	33
Conclusion	36
5. New Estimates of the Uninsured.....	37

Dual Reporting.....	37
Our Approach	40
Revised Estimates of Uninsurance, By Demographic Group.....	40
Discussion.....	42
6. New Estimates of Medi-Cal and Welfare Enrollment Rates.....	45
Pooled Results	45
Discussion.....	47
7. Conclusion	49
Appendix A. Detailed Notes on File Construction and Matching	51
A.1. The Raw Data.....	51
A.2. Matching the CPS and the MEDS.....	52
A.3. Verifying Matches	52
Appendix B. Regression Specification for Response Errors	55
B.1. The Basic Approach	55
B.2 Detailed Logistic Regression Results	57
Bibliography	61

List of Tables

Table 2.1. Medi-Cal Enrollment in California (millions of persons)	6
Table 2.2. Chronology of CPS-ADS Changes and Effect on Health Insurance Coverage	10
Table 2.3. Reporting Rates (CPS/MEDS)	12
Table 3.1. Sample for Matched Analyses.....	14
Table 3.2. Percentage of People from Full Sample in Final (“Matched”) Sample	17
Table 3.3. Congruence in the Matched Sample between MEDS and CPS Data	18
Table 3.4. Mis-Reporting Rates (in %).....	20
Table 3.5. CPS Reference Period	25
Table 3.6. MEDS Data for CPS Imputed Records.....	27
Table 4.1. Adjustment Factors α	35
Table 5.1. Estimates of Dual Coverage and Uninsurance	38
Table 5.2. Health Insurance Coverage Rates: Unadjusted, Adjusted, Discrepancy, Pooled Years	41
Table 5.3. Health Insurance Coverage Rates: Unadjusted, Adjusted, Discrepancy, 2000 Survey Year/1999 Calendar Year	42
Table 6.1. Enrollment Rates: Unadjusted, Adjusted, Discrepancy Medi-Cal, Pooled Years.....	46
Table 6.2. Enrollment Rates: Unadjusted, Adjusted, Discrepancy Welfare, Pooled Years.....	46
Table 6.3. Enrollment Rates: Unadjusted, Adjusted, Discrepancy Medi-Cal, 2000 Survey/1999 Calendar Year	47
Table 6.4. Enrollment Rates: Unadjusted, Adjusted, Discrepancy Welfare, 2000 Survey/1999 Calendar Year	47
Table A.1. Age Differential (CPS Age – MEDS Age) (Percent within Validation Status)	54
Table B.1. Detailed Logistic Regression Results for Medi-Cal	58
Table B.2. Detailed Logistic Regression Results for Welfare	59

List of Figures

Figure S.1 – Reporting Rates (CPS relative to MEDS) by Age and Program	xiii
Figure 2.1 – Medi-Cal Enrollment, in California, by Age, Welfare and Total.....	6
Figure 2.2 – Under-Reporting of Enrollment (CPS/MEDS)	11
Figure 3.2 – CPS Reporting of Welfare Given MEDS Pattern of Receipt.....	23
Figure 5.1 – Dual Coverage Rates and Adjusting Total Health Insurance Coverage (OHI: Other – non-Medi-Cal – Health Insurance).....	38

Summary

High-quality survey data are crucial to our understanding of the effects of the Medi-Cal program in California, and the nation's social welfare system more broadly. We can tabulate the number of people enrolled in Medi-Cal from the official program records, the Medi-Cal Eligibility Data System (MEDS). However, beyond enrollment counts, understanding Medi-Cal's effects often requires survey data because information is needed on both enrollees and non-enrollees. For example, to assess take-up rates we need to know the number of people enrolled as well as the number of people who are eligible for the program. If we want to look at take-up by sub-group, we need more detailed information about the characteristics (e.g., family structure, household income) of enrollees and non-enrollees. If we are interested in assessing overall levels of health insurance coverage, we need information on the full population (enrollees and non-enrollees) and their private health insurance coverage. Addressing policy questions of this form requires survey data.

The Current Population Survey (CPS), Under-Reporting, and Matching

The U.S. Bureau of the Census's March Annual Demographic Survey (ADS) to the Current Population Survey (CPS) is the standard data source for analyses of the Medi-Cal program and the nation's social welfare system more broadly. The CPS is a large (about 50,000 households nationally, 6,000 households in California), household survey with information on program participation (including Medicaid/Medi-Cal and welfare), health insurance coverage, and other household characteristics. Two other features of the CPS data are crucial for policy analyses: (1) The ADS data are collected annually in a relatively consistent manner back to the late 1980s—allowing trend and time series analyses; and (2) The data are released promptly—results of the interviews conducted in March are publicly released in late-August or early-September of the same year—allowing nearly real-time tracking of changes.

Unfortunately, the CPS is known to under-report program participation, including Medi-Cal. The official CPS report notes the problem explicitly:

The Current Population Survey (CPS) underreports medicare [stet] and medicaid [stet] coverage compared with enrollment and participation rates from the Centers for Medicare and Medicaid Services (CMS), formerly the Health Care Financing Administration. A major reason for the lower CPS estimates is that the CPS is not designed to collect health insurance data; instead, it is largely a labor force survey. Consequently, interviewers receive less training on health insurance concepts. Additionally, many people may not be aware that they or their children are covered by a health insurance program if they have not used covered services recently and therefore fail to report coverage. CMS data, on the other hand, represent the actual number of people (who) enrolled or participated in these programs and are a more accurate source of coverage levels.

Furthermore, some analyses suggest that the problem has gotten worse over time.

As we will show below, the under-reporting is substantial, but neither its causes, nor its effects, are well understood. Therefore, with funding from the Medi-Cal Policy Institute and the U.S. Department of Health and Human Services, Administration for Children and Families and the cooperation of the U.S.

Bureau of the Census and the California Department of Health Services (CDHS), we matched individual-level CPS responses to their corresponding MEDS administrative data records. Specifically, as part of its interview, the CPS attempts to collect Social Security Numbers (SSNs) on all respondents age 15 and older. The MEDS data include SSNs for each enrollee. For this project, the Census Bureau supplied a version of the CPS data for 1990 to 2000 that included a scrambled version of the SSN, where available. In addition, the Census Bureau processed a version of the MEDS data for 1989 to 2001 replacing the original SSNs with the same scrambled SSNs. Where possible, we then matched the two files creating a single analysis file with both CPS and MEDS data. To preserve the confidentiality of CPS respondents and Medi-Cal enrollees, the data analysis took place at the UCLA site of the Secure Data Facility of the Census Bureau's California Census Research Data Center. The authors had no access to identifiers (names or Social Security Numbers) and all research results were reviewed to assure that they did not indirectly reveal the identity of or information about CPS respondents or Medi-Cal enrollees.

The Magnitude of Under-Reporting And Our Imputation Model

How serious is the problem of under-reporting? Previous analyses of this question using unmatched data have been limited by the inconsistencies between the two data sources. The CPS, administered in March, asks about program enrollment *at any time* in the last calendar year (i.e., the 2000 CPS asks about program participation between January and December 1999). Aggregate Medi-Cal data is usually reported in terms of persons covered per month. The extent to which discrepancies in aggregate counts based on unmatched data were real as opposed to being an artifact of different data concepts has therefore been unclear. Given the structure of our matched data, we can tabulate the individual level Medi-Cal data from MEDS to be consistent with the CPS questions and thus provide a better estimate of under-reporting in the CPS.

Figure S.1 summarizes that analysis. It considers two age groups (adults—15-65 at the interview, and children—0-14 at the interview) and two program concepts: all Medi-Cal (M) and cash assistance/welfare (W—Welfare). Averaged over the entire period, CPS estimates of total Medi-Cal enrollment for adults are only 70 percent of the counts from the official MEDS administrative data, i.e., Medi-Cal is under-reported by about 30 percent. For children, reporting of Medi-Cal is slightly better, about 75 percent. Unlike some national estimates, there is little evidence of a decline in reporting over time.

This overall pattern in Medi-Cal hides a strong divergence by Medi-Cal sub-program. Enrollment in welfare is severely under-reported. Over the entire time period CPS estimates of total welfare enrollment are only 48 percent of the counts from the official MEDS administrative data. For children, the corresponding figure is 51 percent. For welfare, there is clear evidence of a sharp drop in reporting rates over time. The timing of the drop (in the late 1990s) is nearly simultaneous with the implementation of welfare reform in California (i.e., CalWORKs), perhaps suggesting an increase in the stigma of welfare participation.

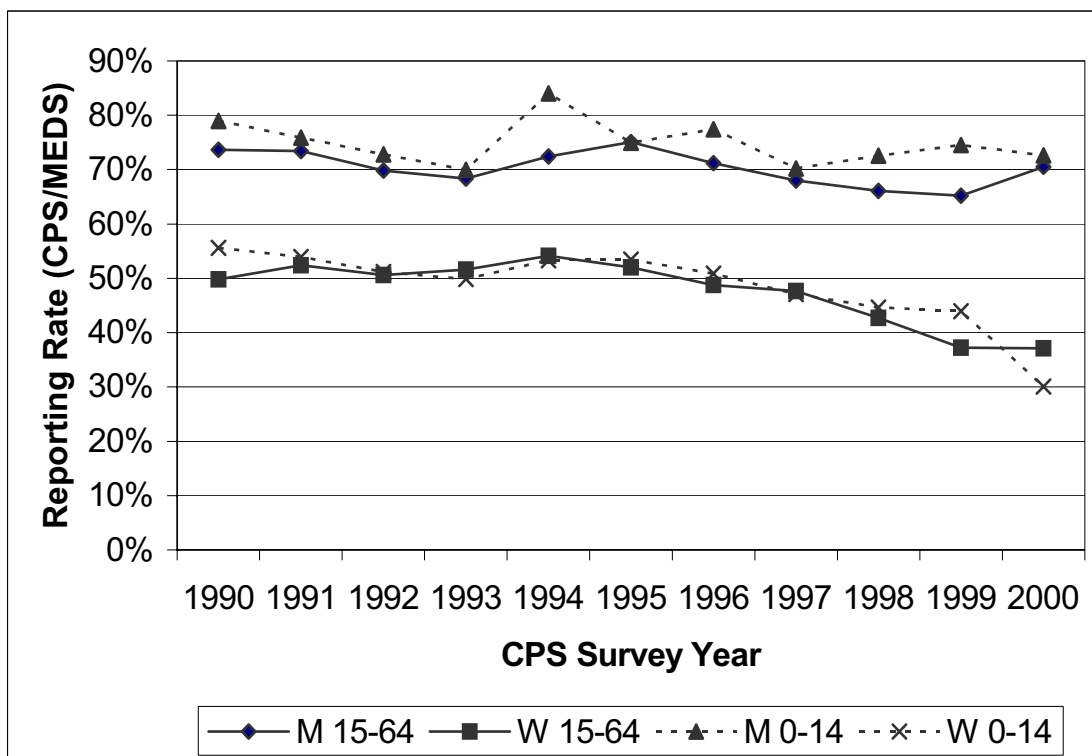


Figure S.1 – Reporting Rates (CPS relative to MEDS) by Age and Program

Source: Tabulations from RAND Merged MEDS File

The under-reporting of program participation in the CPS is severe enough to have substantively important effects on our understanding of the effects of the Medi-Cal program. In this report, we consider two effects. First, under-reporting will lead us to under-estimate take-up rates (the fraction of eligibles enrolled in the program) and thus to over-estimate the need for efforts to increase enrollment or new programs to provide additional coverage. Second, under-reporting will lead us to over-estimate the total number of uninsured people.

Our analysis proceeds as follows. For adults providing an SSN, we overwrite the CPS Medi-Cal responses with the official information from the Medi-Cal administrative data (i.e., treating the MEDS information as the truth). Following Census practice, for children whose parents provide an SSN, we impute Medi-Cal if either the survey response implies Medi-Cal for the child or the administrative data for the parent implies Medi-Cal for the mother (or if not mother, the household head)—in which case the child is almost always covered. However, our ability to match the survey and administrative data is constrained by the fact that only about 62 percent of CPS adults provide an SSN. Furthermore, children under 15 were never asked for an SSN. To address this problem, we build an imputation model to predict mis-reporting among those people without an SSN who we cannot match to the MEDS data. The response errors (i.e., reporting no Medi-Cal in the CPS given actually having Medi-Cal and reporting having Medi-Cal in the CPS given not actually having Medi-Cal) among those not providing an SSN are assumed to follow the general pattern in the sub-sample who do provide an SSN, with an adjustment to force the totals to align exactly (see the full report for details). The problem is more pronounced for children since SSNs are not collected in the CPS for people under age 15. To address this issue we use a combination of information from the head of household and our imputation model. Specifically, where

the mother (or the head of the household if the mother is not present) provides an SSN (as is true for about 66 percent of CPS children), we use the mother or head's Medi-Cal status (from the MEDS or from our imputation) to impute Medi-Cal status to the child. Some Medi-Cal programs include children, but not adults. Therefore, in cases where the child has Medi-Cal, but the head of household does not, the child's data are not changed. Again, as with adults, the imputation includes an adjustment to force the CPS totals (after imputation) to align exactly to the MEDS counts (again, see the full report for details).

These imputations are performed for every observation in the CPS. The resulting individual level file allows us to construct improved estimates of take-up rates and uninsurance coverage. Using the individual-level imputation file, we can consider the effects of under-reporting by respondent characteristics (e.g., gender, age, income).

The Effects of Under-Reporting on Estimates of Medi-Cal Enrollment Rates

If Medi-Cal enrollment is under-reported, then Medi-Cal enrollment rates—the fraction of a demographic group enrolled in Medi-Cal—will also be under-reported. (We note that these are not standard take-up rates that attempt to model actual eligibility from the survey data.) Our analyses of the matched file suggest that the under-reporting is not uniform across sub-groups of the population, so the effects of under-reporting on take-up rates are also not uniform.

Overall Medi-Cal enrollment increases by about 40 percent when we adjust for under-reporting using our imputation model. The increases are slightly larger for adults (42 percent) and slightly smaller for children (38 percent). Consistent with an explanation of under-reporting due to stigma, the increases are smallest for single women with children who are in poverty, largest for those between one and two times poverty, and large for those at more than twice poverty.

Consistent with the even more severe under-reporting, compared to Medi-Cal, the levels of welfare enrollment are lower and the adjustments have a larger effect. The average adjustment more than doubles enrollment rates. The adjustments are similar across children and adults and are smaller for those near poverty, and larger for those out of poverty.

The Effects of Under-Reporting on Estimates of Uninsurance

Another consequence of under-reporting of Medi-Cal enrollment is that it will lead to over-estimates of the rates of uninsurance in the CPS. The magnitude of the over-estimate will depend on the extent to which those under-reporting have other sources of health insurance at some point during the year. If it were the case that everyone who under-reports Medi-Cal did not have any other source of insurance, then we could construct a better estimate of the number of uninsured by subtracting the estimate of under-reporting (i.e., the percent of people in the CPS who report no Medi-Cal, but who our imputation model, based on the matched data, suggests are enrolled) from the raw estimate of the percent of people who are uninsured in the CPS. Conversely, if it were the case that everyone who under-reports Medi-Cal also has private health insurance, then under-reporting would have no effect on the estimates of the uninsured. Our analyses suggest that the truth lies somewhere between these two extremes. Plausibly, we find that under-reporting is more common among those with private health insurance, but under-reporting also includes large numbers of people without private health insurance.

From our matched file, we tabulate rates of other health insurance among people who under-report Medi-Cal. Here we report adjusted estimates of uninsurance based on several different scenarios. We estimate that under-reporting is about 4.1 percent for adults (i.e., 4.1 percent have Medi-Cal but do not report it to the CPS). Consistent with much higher rates of Medi-Cal coverage for children, the corresponding rate of under-reporting is much higher. We estimate that 9.0 percent of all children have Medi-Cal, but do not report it. The unadjusted, or raw, estimate of uninsurance, is 23.5 percent for adults; for children, the unadjusted estimate is slightly lower, 17.8 percent.

The question is: How to use the information from the matched data on the rate of under-reporting to adjust the survey data. A naive approach would, implicitly assume no dual coverage, and subtract the under-reporting from the unadjusted survey estimate of uninsurance. In fact, among those who report in the survey that they have Medi-Cal about a fifth (23.9 percent for adults, 16.7 percent for children) also report private health insurance. Simple tabulations of the matched data suggest that people who have Medi-Cal (according to the survey data), but report that they do not have Medi-Cal in the survey data are much more likely to be dual-covered (i.e., to have other health insurance): 32.3 percent for adults, 34.5 percent for children. Thus, the effect of under-reporting on uninsurance rates is considerable smaller than would be implied by simple subtraction. Using the full imputation model, we estimate uninsurance rates for adults of 20.8 percent (vs. the simple estimate of 23.5 percent) and 11.9 percent for children (vs. 17.8 percent).

Again, adjustments for under-reporting and dual-coverage are not uniform across sub-groups. Sub-groups with higher Medi-Cal receipt rates have larger percentage increases in Medi-Cal coverage with imputation. For adults, the differences across sub-groups are large. For children, the differences across sub-groups are not large.

Summary

We considered the quality of Medi-Cal information in the Current Population Survey, the standard data source for tabulations of Medi-Cal take-up and levels of uninsurance. The analyses are based on an imputation model derived from a match of individual-level survey data with individual-level administrative data for the Medi-Cal program. We find sizable under-reporting of Medi-Cal, leading to sizable under-estimates of Medi-Cal take-up and sizable over-estimates of the fraction of Californians who are uninsured. These results cover the period 1990 to 2000. The Census Bureau made some adjustments to the CPS interview towards the end of this period. Nevertheless, these results suggest caution in basing policy on unadjusted analyses of the CPS data. Analyses based on unadjusted data are likely to substantially overestimate the magnitude of the problem, especially for children.

Acknowledgments

Funding for this analysis was provided by the California HealthCare Foundation (CHCF; through its then separate Medi-Cal Policy Institute—MCPI) and the U.S. Department of Health and Human Services, Administration for Children and Families (DHHS-ACF). The project officers at CHCF/MCPI—Ingrid Aguirre Happoldt, and at DHHS-ACF—Audrey Mirsky-Ashby, Laura Chadwick, and Leonard Sternbach—have been waiting patiently for these results and we appreciate their interest and patience.

This analysis is based on a unique dataset constructed by matching confidential Census Current Population Survey data to confidential administrative data on the Medi-Cal program. Doing so has required the cooperation of several groups. Gene Hiehle and the California Department of Health Services provided the Medi-Cal administrative data and have been supportive throughout this project. B.K. Atrostic and the U.S. Bureau of the Census's Center for Economic Studies provided the Current Population Survey data and handled the matching and deidentification tasks. Senior leadership of the California Census Research Data Center, especially, V. Joseph Hotz at UCLA and Andrew Hildreth at UC Berkeley, have provided crucial support during the negotiations. Center staff, especially Nelson Lim and Becky Acosta, have provided guidance and support with writing the proposal, using the Center, and doing the analyses. The contribution of each of these groups has been necessary in order to gain access to the data.

This research has been presented at RAND, at the California Health Care Foundation, and at the Population Association of America. Comments received in each of these forums have improved the paper. An anonymous RAND interal reviewer also provided very useful comments.

This research is an outgrowth of work begun under the RAND Statewide CalWORKs Evaluation. The constructive comments and guidance provided by CDSS employees during that effort has benefited this effort greatly (though they have not always agreed with our findings). They include Werner Schink, Lois van Beers, Nikki Baumrind, Wilistine Sayas, Aris St. James, and Paul Smilanick.

The data analysis for this project is based on the data preparation work of programmers in RAND's Research Programming Group under the direction of Jan Hanley, who led the effort and did much of the data preparation work herself. Beth Roth, an author on this report, is a member of that group. Other programmers involved in the effort included Christine DeMartini, Laurie McDonald, and Deborah Wesley.

Finally, at RAND this work proceeded within the Labor and Population Program's Center for the Study of Social Welfare Policy. Further information about RAND, the Labor and Population program, and the Center for the Study of Social Welfare Policy can be found at [/www.rand.org](http://www.rand.org), [/www.rand.org/labor](http://www.rand.org/labor), and [/www.rand.org/socialwelfare](http://www.rand.org/socialwelfare). The strong support of RAND, the Labor and Population Program, and its current and former Directors, Arie Kapteyn and Lynn Karoly (respectively), and Assistant Director Rebecca Kilburn, for this effort is gratefully acknowledged. Within RAND, this report has also benefited from the secretarial support of Christopher Dirks and Natasha Kostan. Finally an anonymous reviewer provided technical review improved the final product.

Glossary, List of Symbols, etc.

Symbol	Definition
ADS	CPS March Annual Demographic Survey
AFDC	Aid to Families with Dependent Children
CalWORKs	California Work Opportunities and Responsibility to Kids Act (1997)
CRDC	Census Research Data Center
CCRDC	California Census Research Data Center
CDHS	California Department of Health Services
CDSS	California Department of Social Services
CHCF	California HealthCare Foundation
CHIP	Child Health Insurance Programs (established by statute in 1997, operated by the states)
CMS	Centers for Medicare and Medicaid Services
CPS	Current Population Survey
ESHI	Employer Sponsored Health Insurance
HCFA	Health Care Financing Administration
MEDS	Medi-Cal Eligibility Data System
MCPI	Medi-Cal Policy Institute
NIPA	National Income and Product Accounts
OHI	Other Health Insurance
PIK	Person Identification Key
SIPP	Survey of Income and Program Participation
SSA	Social Security Administration
SSI	Supplemental Security Income
SSN	Social Security Number
TANF	Temporary Assistance to Needy Families

1. Introduction

High-quality survey data are crucial to our understanding of the effects of the nation's social welfare system. If all one wants to know is the number of people participating in a program, then that information can be obtained from administrative data. However, very often, both researchers and policymakers want to know take-up rates (i.e., the fraction of people with certain characteristics enrolled in the program) and the effects of the program on subsequent outcomes (e.g., probability of lacking any health insurance, probability of living in poverty, etc.). For these outcomes, we need richer data that can only be gleaned from surveys; in particular, we need: (1) information on the number and characteristics of nonparticipants; and (2) information on participating families not recorded in administrative data.

Unfortunately, there is considerable evidence that survey data significantly under-report participation in safety-net programs relative to aggregate administrative counts and that the under-reporting has increased over time (e.g, Bavier, 1991). However, most of the evidence to date is based on comparisons between aggregate administrative counts and estimates from survey data. It is our belief that a better understanding of the nature and scope of under-reporting can be obtained by comparing administrative and survey data at the individual level and that is what we seek to do in this report.

This document reports the results of a record-match study of individual-level administrative data for Medi-Cal—the Medicaid program in California, and the Current Population Survey (CPS). With funding from the California HealthCare Foundation (CHCF; through its then separate Medi-Cal Policy Institute—MCPI) and the U.S. Department of Health and Human Services, Administration for Children and Families (DHHS-ACF), and the cooperation of the U.S. Bureau of the Census, the California Department of Health Services (CDHS), and the California Census Research Data Center (CCRDC), we matched administrative data for Medi-Cal from the Medi-Cal Eligibility Data System, (MEDS) to March CPS data for 1990 to 2000. In California, everyone receiving cash assistance (sometimes referred to as welfare)——through Aid to Families with Dependent Children (AFDC), later Temporary Assistance to Needy families (TANF)/California Work Opportunities and Responsibility to Kids (CalWORKs)—is automatically enrolled in Medi-Cal. Since the MEDS administrative data allow us to identify the “type” of Medi-Cal coverage (i.e., why the person is eligible for Medi-Cal), we are able to consider overall Medi-Cal coverage and its two components—welfare and Medi-Cal only (i.e., Medi-Cal, but not welfare)——in our analysis.

Plan of the Report

This report proceeds as follows. The balance of this opening chapter reviews the existing literature on the quality of the CPS data on Medicaid and welfare. The second chapter provides background information on the Medicaid/Medi-Cal program, the MEDS (administrative) data, and the CPS (survey) data. It then characterizes the under-reporting problem, using separate tabulations from each data source. In the third chapter, we turn to the matched data file. For the subset of individuals who provide a valid Social Security Number (SSN), we describe the nature of reporting biases based on a one-to-one match of the survey and administrative data. Unfortunately, not all survey respondents provide an SSN. The fourth chapter provides a technical discussion of our methods for using information from the

matched data to impute welfare and Medi-Cal for the entire California CPS sample. In the fifth chapter, we use the resulting multiply-imputed file to reconsider some of the substantive issues for which the CPS is used. In particular, we explore program take-up by (reported) household income, family structure and other health insurance coverage. The final chapter considers the implications of the results.

Previous Literature on Under-Reporting

The conventional source for information on program take-up is the CPS, the largest annual, national survey. Beginning with the March 1995 CPS, the Census Bureau (Benenfield, 1996a), the Congressional Budget Office (Bilheimer, 1997), General Accounting Office (1997), and the Employee Benefits Research Institute (Fronstin, 1996) each publish annual CPS-based estimates of health insurance coverage and uninsurance. However, the CPS-based estimates of health insurance coverage are much lower and estimates of uninsurance much higher than tabulations from other surveys, such as Survey of Income and Program Participation (SIPP) or the National Survey of America's Families (NSAF) (Bennefield, 1998; Lewis, Ellwood, Czajka, 1998; Fronstin, 2000).¹

In addition and of particular relevance to this study, CPS estimates of Medicaid coverage (Medi-Cal in California) are much lower than corresponding tabulations from administrative data on Medicaid (and Medi-Cal in California), suggesting that survey respondents under-report Medicaid/Medi-Cal coverage. The Urban Institute's TRIM2 model (used by DHHS to simulate program costs) uses administrative data from the Centers for Medicare and Medicaid Services (formerly called the Health Care Financing Administration- CMS/HCFA) to partially correct for such under-reporting. For 1995, this correction for underreporting lowers the fraction of children (0-17) uninsured by 31 percent and the fraction of all non-elderly individuals (0-65) uninsured by 11 percent.²

As part of a discussion of the decline in Medicaid coverage, Ku and Bruen (1999) summarize the national issues and their effect on our understanding of policy.

- 1) "CPS data indicate that about 2.5 million fewer non-elderly people got Medicaid in 1997 than in 1995 (9.3 percent fewer), while administrative data indicate that 1.2 million (3.2 percent) lost Medicaid."
- 2) "CPS data indicate that more children lost coverage than adults from 1995 to 1997, while administrative data indicate [that] the declines were larger for adults."
- 3) "[T]he total number of nonelderly people who had Medicaid at any time in a given year was about 25 to 30 percent lower in the CPS than in administrative counts."
- 4) "[T]here appears to be a growing discrepancy between CPS and administrative data concerning the receipt of benefits like Medicaid, welfare, and food stamps in recent years. . . . Using measures of enrollment during the year, the CPS Medicaid

¹ Other papers focusing on question wording for health insurance items include Rajan et al. (2000), and Nelson and Mills (2001).

² For similar comments about welfare, see <http://www.census.gov/hhes/www/income/assess1.pdf>.

participation estimates were 75 percent of administrative counts in 1995, but fell to 70 percent in 1997.”³

- 5) “Some believe that respondents to the CPS may be reporting their current insurance status, rather than answering the actual question about insurance at any time in the prior year.”

Such reporting biases would cause over-estimates of the number of uninsured Americans and, thus, of the demand for the programs being created by new policy initiatives. Lower than expected enrollment has in fact been a problem (Alpha Center, 2000). While simple reporting bias is unlikely to explain all the lower than expected enrollment, such reporting bias has explicitly been cited by some observers (e.g., Alpha Center, 2000).

The problem of under-reporting is perceived to be so severe that the official U.S. Bureau of the Census report on health Insurance (P60-220, 2002) notes it explicitly and at length in its “Technical Note”:

The Current Population Survey (CPS) underreports medicare [stet] and medicaid [stet] coverage compared with enrollment and participation rates from the Centers for Medicare and Medicaid Services (CMS), formerly the Health Care Financing Administration. A major reason for the lower CPS estimates is that the CPS is not designed to collect health insurance data; instead, it is largely a labor force survey. Consequently, interviewers receive less training on health insurance concepts. Additionally, many people may not be aware that they or their children are covered by a health insurance program if they have not used covered services recently and therefore fail to report coverage. CMS data, on the other hand, represent the actual number of people (who) enrolled or participated in these programs and are a more accurate source of coverage levels.

The problem of under-reporting appears to be particularly severe for welfare. Welfare recipients are categorically eligible for Medicaid. In fact, the CPS imputes Medicaid to anyone who reports receiving welfare. However, welfare reform appears to have worsened reporting of welfare in the CPS, perhaps because of confusion over program names, perhaps because of increased stigma of welfare receipt.

³ Furthermore, concern about the problem has increased. See, for example, Levit et al. (1992, pp. 45-46), reflecting minimal concern about undercounting. “CPS counts of people covered by Medicare and Medicaid programs are reasonably consistent with Health Care Financing Administration (HCFA) program data after allowing for the institutional component missing from CPS.” They compare the 1991 CPS estimate (for 1990) of 24.3 million persons to the HCFA Medicaid program estimate of 25.3 persons. They attribute the difference (only about 4 percent) to “the institutionalized population not included in CPS and difficulties that surveys have capturing Medicaid recipients.” They note that estimates of change over time (in particular 1980 to 1991) are quite similar across the CPS and HCFA data.

See Fronstin (1997), HCFA (1996), and Lewis, Ellwood, and Czjaka (1998) for claims that Medicaid under-reporting has increased.

2. Medi-Cal, the CPS, The MEDS, and Under-Reporting

The core of this project is a data match between administrative data for California’s Medicaid program—Medi-Cal (i.e., the MEDS data)—and CPS data. This section begins with a brief description of the Medi-Cal program. It then describes the administrative data (the MEDS) and the survey data (the CPS). Finally, we provide some simple tabulations using the *unmatched* data.

The Medi-Cal Program

Since 1965, Medicaid—a joint federal-state program—has provided health insurance to current welfare recipients and some other qualifying families. During the 1980s and 1990s, coverage was significantly expanded, with particular attention to poor children (often referred to as “the percent programs”) and families that are welfare-eligible, whether they are actually on welfare or not (the 1931(b) program).⁴ California’s Medicaid program, Medi-Cal, is a joint effort of the California Department of Health Services (CDHS), which administers the program and handles payments, and the California Department of Social Services (CDSS), which supervises county welfare departments that handle enrollment and re-enrollment.⁵

Figure 2.1 (and Table 2.1) shows Medi-Cal enrollment from the MEDS data (described below) according to the CPS concepts we will use in our main analysis. In particular, we tabulate the total number of individuals enrolled at any time in the calendar year. We distinguish welfare from other Medi-Cal (Medi-Cal only, or simply “MO”). Finally, we consider only the non-elderly, in two groups: those 0-14 as of the following March (who we refer to as “Children”) and those 15-65 as of the following March (who we refer to as “Adults”; we discuss the reason for this child/adult break at 14/15 below).

In the late 1980s, Medi-Cal had just over 4 million enrollees, approximately evenly divided between adults and children, with more welfare than Medi-Cal Only (i.e., during the calendar year at least some Medi-Cal, but never welfare). During the early 1990s, the number of enrollees grew rapidly to over 6 million because of a combination of two factors. First, program eligibility was deliberately expanded. Second, California’s deep recession made more people income-eligible, especially through rapid growth in welfare/cash assistance.

From the mid-1990s to the early 2000s, Medi-Cal enrollment has remained relatively stable, near 6 million. This stability is the result of offsetting trends in Medi-Cal sub-programs such as welfare and 1931(b). First, as in the rest of the nation, there has been a sharp drop in welfare/cash assistance since the early 1990s. Second, as was intended (but after a transition period), the 1931(b) program’s growth has more than offset the shrinkage in cash assistance. Third, the other components of Medi-Cal such as Supplemental Security Income (SSI), Medically Needy, and “Other” are relatively stable.

⁴ For more discussion of these eligibility changes and their effects, see Gruber (2000).

⁵ For more information on Medi-Cal and its multiple programs see: <http://www.medi-cal.org/> and its fact sheet: <http://www.medi-cal.org/resources/view.cfm?section=Resources&itemID=1397>. For more information on the administration of Medi-Cal, see Klerman and Cox (2003).

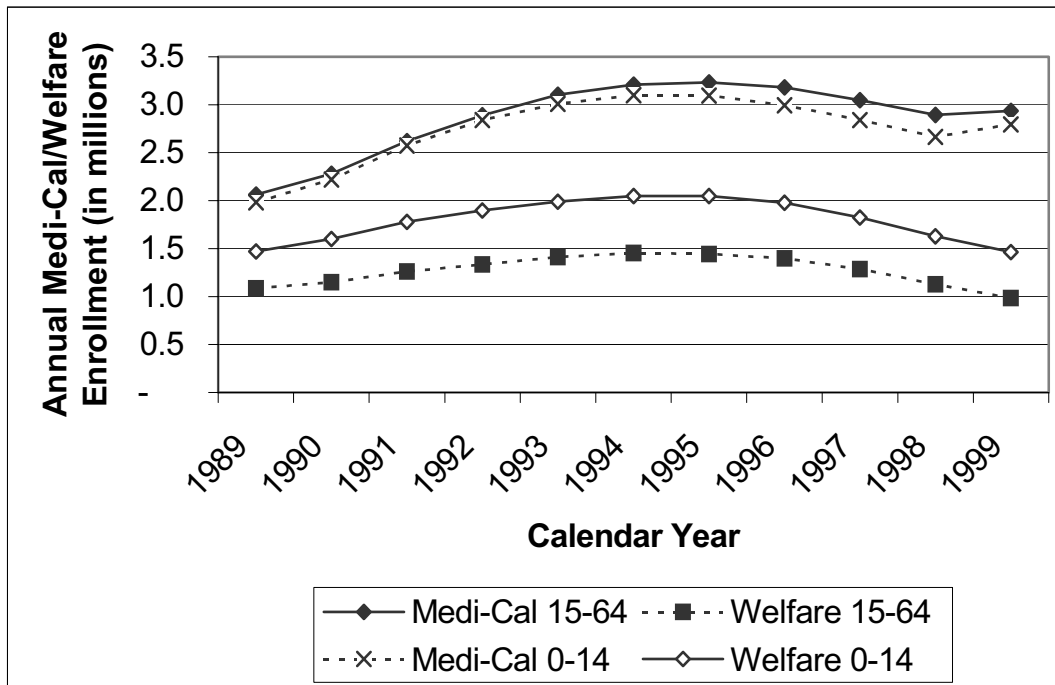


Figure 2.1—Medi-Cal Enrollment, in California, by Age, Welfare and Total
 Source: Tabulations from RAND Merged MEDS File

Table 2.1.
 Medi-Cal Enrollment in California (millions of persons)

Year	Adults			Children		
	M	W	MO	M	W	MO
1989	2.1	1.1	1.0	2.0	1.5	0.5
1990	2.3	1.2	1.1	2.2	1.6	0.6
1991	2.6	1.3	1.4	2.6	1.8	0.8
1992	2.9	1.3	1.6	2.8	1.9	0.9
1993	3.1	1.4	1.7	3.0	2.0	1.0
1994	3.2	1.5	1.8	3.1	2.1	1.1
1995	3.2	1.4	1.8	3.1	2.0	1.0
1996	3.2	1.4	1.8	3.0	2.0	1.0
1997	3.0	1.3	1.8	2.8	1.8	1.0
1998	2.9	1.1	1.8	2.7	1.6	1.0
1999	2.9	1.0	2.0	2.8	1.5	1.3

Source: RAND tabulations from merged MEDS file.

Note: M-Medi-Cal, W-Welfare, MO-Medi-Cal Only. Tabulated according to CPS concepts: Enrollment is at any time in calendar year (not in each month); Adults are 15-65 in March of the next year; Children are 0-14 in March of the next year.

The Medi-Cal Eligibility Data System (MEDS)

Real-time enrollment information in Medi-Cal is maintained in the MEDS. County welfare departments update this system as individuals are enrolled in or drop out from the program. Providers

check the system to verify whether an individual is covered or not, and which services would be reimbursed by the Medi-Cal program.

Specifically, the MEDS is a person level file using Social Security Number (SSN) as its primary person-level identifier. The data comes to us in overlapping 15-month batches. We link the records using SSN. The quality of the SSN information is expected to be relatively high because the MEDS program has an aggressive program to catch and correct errors in SSN. Moreover, SSN is a primary variable for fraud detection. It is verified initially using the Social Security card and is regularly re-verified. There are some issues with pseudo-SSNs assigned to children prior to their receiving an SSN, but these are also addressed in large part by the state program to identify and fix problem records.

For each month, the file contains an “aid code”. This aid code indicates the program through which the individual receives Medicaid/Medi-Cal coverage (e.g., welfare, SSI, 1931(b), Medically Needy). These codes are of considerable operational relevance. They determine for what medical services a provider will be reimbursed and whether or not a co-pay (referred to as a “share of cost”) must be collected.

We code anyone with a non-zero aid code as receiving Medicaid/Medi-Cal. In the MEDS, the relevant concept is coverage; i.e., is this individual currently holding a Medi-Cal card such that a provider could deliver services with a reasonable expectation of reimbursement. We believe that this corresponds relatively closely to the concept about which the CPS attempts to collect a response. It is possible (in fact, given some ways in which eligibility is conferred and continued, it is likely) that some people are enrolled (or remain enrolled) in Medicaid/Medi-Cal without realizing it. Since we can not directly identify such individuals, we have no choice but to include them. We note that for many (but not all) purposes, they should be included (e.g., if they would be admitted to a hospital for emergency care, the hospital would be reimbursed). Throughout the balance of this document, we use the terms “covered”, “enrolled”, and “participating” interchangeably. In particular, we note that while some authors use “participating” to those who actually receive services, we intend no such distinction.

In California both before and after welfare reform, cash assistance (AFDC/CalWORKs) automatically confers eligibility for (i.e., coverage by) Medicaid/Medi-Cal with the highest level of covered services and no share of cost. As noted above, for operational reasons, this source of coverage is noted in the Aid Code on the MEDS file. It is conventionally used by the state for official purposes (e.g., selecting the Quality Control sample) and by outside analysts to explore welfare related issues. We code welfare corresponding to any AFDC/CalWORKs code providing cash assistance. This includes one-parent and two-parent cases and California’s special refugees and legal immigrant programs (some of which are “state only”). They are referred to as “CalWORKs” in the program names and program documents. In almost all cases, program rules explicitly attempt to treat individuals in all such sub-programs identically. Thus, it seems plausible that CPS respondents in each of these sub-programs would also answer that they were in AFDC/CalWORKs.

Individual-level extracts from the file provide a complete historical record of Medi-Cal eligibility for 1987 forward. Crucially for our purposes, the file includes linking information (name, Social Security Number—SSN), some basic demographics (gender, date of birth, race/ethnicity), and detailed Medi-Cal program information.

In our analysis below, we treat the MEDS records as “truth.” This is a reasonable approximation given their use by providers in determining whether care will be reimbursed. However, the MEDS data are not always absolutely correct. Careful study of the MEDS data suggests some anomalies when

counties had trouble updating the records (e.g., for two months in late 1990, there is a period of a few months when there appear to be no entries onto welfare for Los Angeles County). Card, Hildreth, and Shore-Sheppard note some seam bias (sharp increases in transition rates across versions of the MEDS file), which also suggests some reporting error in the MEDS.⁶

In addition, we note that some people may be enrolled in Medi-Cal but might not be aware of it. In particular, the *Edwards v. Kizer* decision requires California's counties to continue Medi-Cal eligibility for welfare leavers until their eligibility for continued Medi-Cal can be determined. Moreover, California's implementation of the Medicaid 1931(b) program and the provisions of California SB 87 have the effect of keeping many welfare leavers on Medi-Cal even without filing a new application.⁷ It is widely believed that many of these people do not realize they are covered.

The Current Population Survey (CPS)

The CPS is a monthly survey of about 50,000 households conducted by the U.S. Bureau of the Census for the U.S. Department of Labor.⁸ The CPS's primary purpose is to provide official monthly estimates of the unemployment rate, a key business cycle indicator. With its associated sampling weights, it is representative of the American non-institutional population.⁹

Since 1948, in its spring survey the CPS has included additional questions on annual income in the previous year.¹⁰ Today, those additional questions are asked at the end of the March survey (corresponding to the arrival of W-2s and household preparation of federal income tax returns) and are referred to collectively as the Annual Demographic Survey (ADS). Over the years, the set of supplementary questions has grown.

Most important for our purposes, since 1980, the ADS has included detailed questions on health insurance coverage and welfare receipt in the previous calendar year (not as of the date of the March

⁶ See Card, Hildreth, and Shore-Sheppard (2001) for some further discussion of these issues. The seam bias problem should be less severe in the annual reference period of the CPS which we analyze than in the monthly reference period of the SIPP that Card, Hildreth, and Shore-Sheppard analyze. Note also that their biggest matching problems are with children, for whom we do not have SSNs and therefore do not match. Finally, note that below we limit our sample to the validated records which should increase the quality of the SSN data.

⁷ Medicaid Section 1931(b) was a new program created by federal welfare reform (the Personal Responsibility and Work Opportunities Act of 1996) to guarantee Medicaid to any family that would have been eligible for welfare before welfare reform. Section 1931(b) also gave states the option of expanding 1931(b) eligibility to align it with eligibility for cash assistance. California did so with the net effect that welfare leavers with income up to about 165 percent of the poverty line remain indefinitely eligible for Medi-Cal. In practice, implementation of Section 1931(b) in California was delayed until early 1999, but indirect effects (the "Edwards Hold," see Klerman and Cox, 2004) were felt beginning in early 1998.

California SB 87 (chaptered September 30, 2000, effective July 1, 2001) streamlined continued enrollment in Medi-Cal for welfare leavers through adoption of an ex parte process and, in practice, a presumption of continued eligibility for Medi-Cal among welfare leavers. This implementation occurred after the period covered by our data.

⁸ For more on the CPS, see <http://www.bls.census.gov/cps/overmain.htm>.

⁹ The restriction of the CPS universe to the non-institutional population is potentially problematic for analyses of Medi-Cal. While most Medi-Cal enrollees are young, most Medi-Cal expenditures go to the elderly in nursing homes. That group is not in CPS's universe, which is the non-institutional population. The MEDS data do not have a flag for institutional residence. As a partial correction, our analyses below exclude those age 65 and over.

¹⁰ For more on the March Annual Demographic Supplement to the CPS, see:

<http://www.bls.census.gov/cps/ads/adsdes.htm>.

interview).¹¹ These questions began as an attempt to expand the definition of “income” to include employee benefits and non-cash government benefits (Food Stamps, subsidized housing, medical assistance, etc.).

Specifically, in the battery of questions on income receipt, the CPS asks adults about the receipt of public assistance (PAW_YN) by the household in the previous year. Those who respond in the affirmative are then asked the specific form of public assistance (PAW_TYP). We exclude from our welfare variable those who report public assistance, but not AFDC/CalWORKs. This group would include those receiving Supplemental Security Income and those receiving General Assistance (i.e., county level assistance to families without children). This definition corresponds to what is recorded in the MEDS administrative data; i.e., anyone receiving Medi-Cal due to California’s welfare program (called CalWORKs post-welfare reform). See below for a discussion of how we handle imputed and allocated responses (in each analysis).

The definition is, however, problematic for households where some people are on welfare and others are not. In this case, we have followed standard CPS analysis practice and imputed welfare to everyone in the household (more precisely the family or sub-family) when anyone in the household reports welfare. This is likely to lead to false-positives. If someone in the household receives welfare, but this individual does not; we will incorrectly assign them welfare in our CPS analysis file. This is a serious problem with the way the CPS collects its data. It appears unremediable given the available CPS data. We note that from other analyses, it appears that this problem is most salient for child only cases where the parent is an undocumented immigrant (and thus ineligible for welfare), but the child was born in the U.S. (and thus a citizen, eligible for welfare).

Combining the questions on health insurance as an employee benefit with the questions on participation in government health insurance programs yielded a rough measure of total health insurance coverage; its complement provided an estimate of those without health insurance. Until the 2000 interview, the last included in our analysis file, there was no direct question in the CPS about being uninsured. Rather uninsured status is inferred from negative answers to questions about receipt of Medicaid and other types of health insurance (see below Table 2.2 for a summary of changes to the CPS questionnaire).

Specifically, we assign Medicaid/Medi-Cal based on the Census Bureau’s composite Medicaid variable. This variable combines information from the initial CPS question about Medicaid (MCAID—“Was ... covered by medicad) and also a set of follow-up questions among those who otherwise report no other health insurance (OTHYPn). Equivalently, uninsurance is coded based on negative answers to the specific probes for each type of health insurance and to the follow-up questions. See below for a discussion of how we handle imputed and allocated responses (in each analysis).

As this discussion suggests, the individual questions were not originally intended to generate an estimate of the size of the population without health insurance. With issues of uninsurance becoming more salient, in 1988, the Census Bureau refined the questions.¹² Questions about employer-based

¹¹ This discussion draws on Nelson and Mills (2001).

¹² In addition, in 1983, the Census Bureau began a second national survey, the SIPP). The SIPP is a moderate-sized panel survey with more detailed questions on income and program participation (as its name implies). The original vision appears to have been that the SIPP would replace the CPS-ADS for many purposes, including the measurement of health insurance. However, for a variety of reasons (including varying sample size, issues related to

health insurance that previously had only been asked of employed individuals were asked of all individuals 15 or older, regardless of whether they worked. This change should have captured retiree coverage and COBRA benefits (i.e., benefits from a previous employer). In addition, for children, questions were added about health insurance coverage from individuals not residing in the household. This change should have captured coverage provided by non-coresident parents. Finally, the imputation methods for children's coverage were revised and additional questions on Medicaid were added (see Levit et al., 1992; Moyer, 1989; Swartz and Purcell, 1989; and EBRI, 2000).

Additional changes have been made since then. (See also Swartz, 1997; EBRI, 2000.) Table 2.2 presents a detailed chronology. Census analyses suggest that the changes in survey years 1996 and 2000 increased reported health insurance coverage by about 1 percentage point each.

Table 2.2.
Chronology of CPS-ADS Changes and Effect on Health Insurance Coverage

Year	Change	Effect on Health Insurance Coverage
1981	First health insurance questions (employer sponsored and government sponsored) on CPS-ADS	<Not applicable>
1988	Introduction of new CPS processing system	Minimal
1989	Addition of questions on child health insurance coverage (previously coverage of children was imputed based on adult responses)	Moderate
1993	Switch to 1990 Census population controls	Minimal
1994	Switch from paper and pencil to Computer Aided Personal Interviews	Minimal
1996	Questions reordered and modified to improve information on Medicaid	Possibly moderate; see Swartz (1997).
1998	Indian Health Service no longer considered coverage	Minimal
2000	Verification questions added	Moderate (about 1 percentage point)
2001	Switch to 2000 Census population controls	Minimal (less than 1 percentage point)
2001	Addition of questions on state CHIP programs	
2002	Additional sample (78,000 rather than 50,000) to estimate state health insurance coverage rates	Minimal (less than 1 percentage point)

Note: "Year" refers to the survey year. The CPS questions refer to the previous calendar year; i.e. "2001" refers to the survey conducted in March of 2001, collecting information about calendar year 2000. Effect is on the total national coverage rates. Purely because of sampling issues, effects are larger at the state level. Because of substantive issues, effects are often larger for components (e.g., Medicaid). "Minimal" is less than 1 percentage point; "Moderate" is more than one percentage point.

Corresponding to the fact that California has 12 percent of the national population, the annual March ADS to the CPS has about 6,000 California households, about 13,000 individuals (adults and children), and about 2,000 individuals on Medi-Cal.

its panel structure, and slow data release; see Short, 2001), the CPS remains the primary data source for counts of the uninsured. For an analysis similar to this one for the SIPP, see Card et al. (2001).

Levels and Trends In Mis-Reporting

Before turning to the more detailed results from the matched data, we conclude this section with an analysis comparing simple (unmatched) tabulations from the two data sources (i.e., a comparison of the aggregate administrative counts and the CPS estimates of the population enrolled). Figure 2.2 plots the ratio of CPS enrollment to MEDS enrollment. (Table 2.3 provides the underlying numbers.) These tabulations are made using the individual-level MEDS files. We have aligned the counts to match the CPS concepts so that “enrollment” is defined as being recorded in the MEDS as enrolled at any time in the past year. Age is as of March of the next year. The division between adults and children follows the CPS at 14/15 at the interview. California residence in the CPS is ascribed based on residence a year before the interview (not at the interview, as is conventionally done).

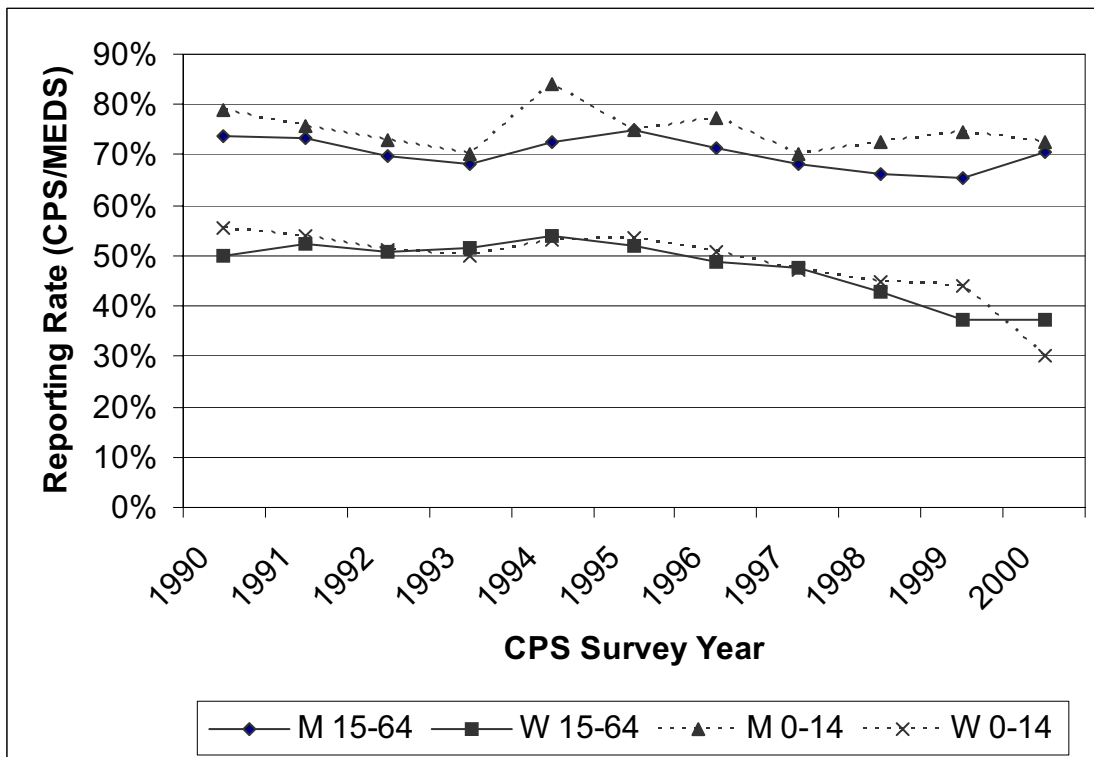


Figure 2.2 – Under-Reporting of Enrollment (CPS/MEDS)
 Source: Tabulations from RAND Merged CPS File and RAND MEDS file

Table 2.3.
Reporting Rates (CPS/MEDS)

	Adults		Children	
	Medi-Cal	Welfare	Medi-Cal	Welfare
1991	73%	52%	76%	54%
1992	70%	51%	73%	51%
1993	68%	52%	70%	50%
1994	72%	54%	84%	53%
1995	75%	52%	75%	53%
1996	71%	49%	77%	51%
1997	68%	48%	70%	47%
1998	66%	43%	73%	45%
1999	65%	37%	75%	44%
2000	71%	37%	73%	30%

Source: RAND tabulations from RAND CPS matched file and the RAND MEDS file.

Note: M-Medi-Cal, W-Welfare. Tabulated according to CPS concepts: Data refers to calendar year preceding the survey year. Enrollment is at any time in calendar year (not in each month); Adults are 15-65 in March of the next year; Children are 0-14 in March of the next year.

As the figure and the table show, there clearly is under-reporting. For both adults and children, CPS Medi-Cal counts are only about 70 percent of MEDS Medi-Cal counts. CPS welfare counts are only about 40 percent of MEDS welfare counts.

Unlike the earlier characterization of the national data, there is little evidence of trend in California's Medi-Cal reporting rates. Perhaps there is some increase in Medi-Cal for adults in 2000 with the latest changes. There is a slight increase for children in 1996. The situation is very different for welfare. Welfare reporting rates for adults have fallen from about 50 percent to about 40 percent over these eleven years. For children, there is a large additional drop in 2000.

Discussion

This section has provided basic background information. It described the Medi-Cal program and the two data sources. It then compared the two data sources in aggregate to provide a rough characterization of the mis-reporting. In the next section, we turn to the match: what it can contribute and what are the technical issues in using the information provided.

3. The Matched Data

In this chapter, we turn to the matched data so that we can investigate mis-reporting at the individual level. In order to understand mis-reporting, we begin by creating and analyzing a dataset containing only the highest-quality matches. In this chapter, we describe the congruence of the CPS responses to the MEDS information, where we treat the MEDS information as “truth.”

What we are able to report is strongly constrained by Census disclosure rules. To preserve the confidentiality of CPS respondents, those rules limit our ability to report exact results for tables with small cell sizes in any of the cells. Thus, to preserve confidentiality, for most analyses, we pool observations across all of the survey years. In some cases, where the cell sizes are particularly small, we combine cells or report the predictions from simple regression models.

As we discuss here and in the next chapter, not everyone provides an SSN, so the matched data cannot provide a complete characterization of the quality of reports. We defer until the next chapter a discussion of how we use the information from the matched sample to make inferences about the entire population—both those who do provide an SSN and those who do not.

Matching

Our first task in constructing the matched dataset is to match the administrative data (from the MEDS) to the survey data for California residents (from the CPS) where possible. Table 3.1 summarizes the sample selection rules for the analyses in this chapter. Appendix A describes the matching process, sample selection rules, and results in more detail. We note that we deliberately suppress the exact sample counts in each year to preserve our ability to present other more substantive results by year later in the analysis.

Our primary focus is on “adults” (defined as age 15 to 65 as of the March interview) residing in California. The reason for these two exclusions is straightforward. First, the CPS only collected SSNs for those aged 15 and older, so there is no possibility of matching “children” (defined as those under age 15). Second, the MEDS is administrative data from the state of California and, thus, only includes information on people enrolled in California’s Medicaid program, Medi-Cal.

Table 3.1.
Sample for Matched Analyses

Sample	Validated	Unvalidated	Total
Adults 15-65			
Initial	53,515	46,109	99,624
No SSN	18,677	18,697	37,374
Movers	1,169	1,105	2,274
Imputed Data	1,519	1,007	2,526
Bad Match	446	704	1,150
Final	31,704	24,596	56,300
Children, 0-14			
Initial	21,668	18,925	40,593
No Parent	251	242	493
No Parental SSN	6,910	6,537	13,447
Parent Mover	415	366	781
Parent Imputed Data	568	432	1,000
Parent Bad Match	237	426	663
Final	13,287	10,922	24,209

Note: Validated Years—1991, 1994, 1997-2000; Unvalidated Years—1990, 1992, 1993, 1995, 1996.

The basic sample begins with California adults. Consistent with the CPS questions on program participation that refer to participation in the previous year, we define “California adults” based on state of residence a year prior to the interview.¹³ From this sample of California adults, the “matched sample” drops the following groups:

- *No SSN*: Even for those age 15 and older, not everyone has a (scrambled) SSN on our internal CPS file.¹⁴
- *Movers*: Our MEDS data only include welfare receipt in California. So for movers, responses about Medicaid and welfare might refer to enrollment in a different state, which would not be

¹³ Using state of residence a year before the survey is more consistent with the CPS reference period than using state of residence at the survey. We note, however, that it does not appear to be standard practice in analyses of program participation or health insurance coverage.

We also note that the CPS reference period induces some standard coverage issues. The sample is drawn in March of the following year and is, therefore, not completely representative of everyone alive in the reference year (or even as of the end of the reference year). The divergence will include births and deaths and changes in residence (in the United States at all, in an institution).

¹⁴ In 1990, 1992, 1993, 1995, and 1996, our file includes “unvalidated SSNs,” i.e., the SSNs are simply a scrambled version of the SSN provided by the CPS respondent.

In 1991, 1994, and 1997-2000, our file includes a validated version of the SSN, not (a scrambled version of) the SSN provided by the CPS respondent. Exact details on the validation process are not available, but it appears that Census provided a list of SSNs and basic demographic information (name, gender, age) to the Social Security Administration (SSA). SSA then cross-checked the information against its SSN records. For some cases, an SSN was imputed onto [into?] the file based on name, place of birth, and birth date; for other cases, a provided SSN was deleted based on failure to match on these criteria.

Some of the tabulations below (including Table 3.1) tabulate results separately by validation status of the survey file (i.e., the entire year). To avoid ascribing lack of congruence between CPS and MEDS data to improper SSNs, other analyses below drop all unvalidated years.

recorded in the MEDS data. For the analysis in this chapter, we therefore also drop movers.¹⁵ (See further discussion of this issue at the end of Appendix A.)

- *Imputed Data:* Imputed data on enrollment are not informative for the quality of unimputed responses. Imputed data on the matching variables might cause us to incorrectly accept or reject an SSN match. (See the next bullet.) We therefore drop from the analyses in this chapter anyone with imputed responses on program enrollment (Medicaid or welfare) and anyone with imputed responses on the basic demographics used in matching (gender and age; there are only a trivial number of such individuals).
- *Bad Match:* To accept a match, we require a match on gender and on age plus or minus one year.

Note that even though we refer to this as our “matched sample,” it is more properly the sample of people who *could have* matched to the MEDS (i.e., they provided an SSN). We do *not* require an actual match to the MEDS because we want to include people in our sample who did not receive Medi-Cal (who may or may not report their program participation correctly in the CPS). Our MEDS extract contains a record for each individual who has received Medi-Cal during 1989 to 2001. Individuals who did not receive Medi-Cal during this period should not and will not appear in the MEDS data. We leave them in the matched sample and infer that they never received Medi-Cal. We note, however, that for this group it is not possible to verify that the gender and age match across the two data sources (i.e., if we had successfully matched based on SSN, we might have found that the gender and age did not match across the MEDS and the CPS).

Finally, the analyses in this chapter consider only the “validated” data in order to generate an analysis file with only the highest quality matches. When the CPS collects SSNs, it also asks permission to “validate” SSNs with the Social Security Administration (SSA). We were unable to obtain a formal written description of how validation was done for these data files. It appears that for a respondent who does not provide an SSN, Census passes his/her name, gender, and age to SSA who attempts to impute an SSN from the official SSAN SSN files (the NUMIDENT file). In addition, it appears that the SSNs passed to SSA are checked and either replaced or deleted. The presumption is that this validation improves the quality of the matches. In some years (1991, 1994, and 1997-2000), Census provided us with the validated files; in other years (1990, 1992, 1993, 1995, 1996), Census provided us with only the unvalidated files. Presumably, validation improves the quality of the SSNs. Our results below support that presumption.

The CPS’s failure to collect SSNs for children implies that we cannot directly apply similar matching methods to children. In the substantive analyses below, we make a rough imputation of the implications of our analysis of the matched data for children and, therefore, for the entire population. To do so, we impute Medi-Cal based on parental Medi-Cal receipt.

There are two issues in doing so: (1) Which parent? And (2) How to impute? As to which parent, we attempt to use the mother. Specifically, we use the CPS PARENT variable. If that person is female, we use her Medi-Cal information. If that person is male, we use the Medi-Cal information for that person’s spouse. If this algorithm does not identify a female, we use the CPS family’s reference person.

¹⁵ Note that our initial sample was defined as people in California a year before the survey. Therefore, the sample we drop as “movers” are those who were in California a year before the survey, but not in California at the survey.

Having identified the person, we impute Medi-Cal to the child if *either* the parent or the child had Medi-Cal, where the imputation is based on the MEDS information if we matched the parent to the MEDS data or our multiple imputation if the parent did not match to the MEDS data. Table 3.1 also gives the total number of California children and the number of them who have a reference person in the matched sample.

The net effect of the sample restrictions outlined above is that although the CPS has 99,624 California adults during the period 1990 to 2000, our narrow matched sample for data quality analysis is only a third of that, 31,704. Only six of the eleven years are validated, cutting our sample nearly in half. About a third of the validated sample provides an SSN, about 2 percent of the sample in California a year before the survey is not in California at the survey, about 3 percent of the sample has imputed data, and about 1 percent is a “bad match” (where the counts are hierarchical; applying each criteria in turn).

For our analyses of CPS data quality, we will use these narrowly defined/highest quality matches. For our extrapolations to the full file, we will use all the matches, making multivariate corrections for the effects of validation (see the discussion of details in the next chapter). As expected, we note that validation increases the fraction of adults with SSNs (from 59 percent to 65 percent). However, the bad match rate nearly doubles with validation (from 0.8 percent to 1.5 percent). Apparently, some of the SSNs added at validation are incorrect.

In the verified sample, we tabulate the weighted fraction of people with given characteristics in the final (i.e., matched) sample (see Table 3.2). This is a rough proxy for presence of an SSN; about 86 percent of those in the full sample, but not in the final sample, are dropped because of a missing SSN. Overall, 58 percent of adults are in the final sample; i.e., provide an SSN, are not movers, do not have imputed data, and are not “bad matches,” (These figures are weighted, unlike those in Table 3.1, which are unweighted).

Based on our descriptive analysis, providing an SSN does not appear to be random. More advantaged people are more likely to provide an SSN and thus to appear in the final sample; less advantaged people are less likely to provide an SSN and be in the final sample. Relative to the sample overall (58 percent), minorities (black and Hispanic, 50 percent), high school drop-outs (50 percent), and those in poverty (less than half the poverty line (45 percent), half the poverty line to the poverty line (50 percent), the poverty line to one and a half times the poverty line (51 percent), and one and a half times the poverty line to twice the poverty line (52 percent)), and single females with children (54 percent) are less likely to be in the final sample. Those with at least some college (63 percent) and other (non-Medi-Cal) health insurance (63 percent) are more likely to provide an SSN.

Given that they are less socially advantaged, we might expect those enrolled in Medi-Cal to be less likely to provide an SSN. Offsetting this, Medi-Cal enrollees are required to supply an SSN each time they deal with the welfare office, so they are likely to know their SSN, and perhaps be less reluctant to give it out. In fact, those reporting to the CPS that they have welfare are more likely than others to provide an SSN (65 percent); those reporting to the CPS that they are enrolled in Medi-Cal are as likely as those not reporting Medi-Cal to provide an SSN (58 percent).

Combining both the validated and unvalidated years, our basic sample has 40,593 children. For almost all of them, we can locate a parent on the file. For about 66 percent of those parents, we have an SSN. This rate is similar to the 62 percent of adults that provide an SSN (again, pooling the validated and unvalidated years). As with the adults, we do not use information from the administrative data in the

case of a move, imputed data, or a bad match. This comprises about six percent of all children; the corresponding figure for adults is also six percent.

Table 3.2.
Percentage of People from Full Sample
in Final (“Matched”) Sample

Sample	% in Final Sample
Overall	58%
Male	58%
Hispanic	50%
Black	50%
HS Drop Out	50%
Some College	63%
FPL<0.5	45%
0.5<FPL<1.0	50%
1.0<FPP<1.5	51%
1.5<FPP<2.0	52%
Kids in Household	57%
Single Female w/Kids	54%
Other Health Insurance	63%
Welfare	65%
Medi-Cal	58%

Source: Weighted tabulations from RAND CPS-MEDS Match file.

Note: FPL—Federal Poverty Line.
Exclusions are primarily no SSN (see Table 3.1)

Simple Reporting Rates

We have deliberately constructed this sample to maximize the congruence between the survey data and the administrative data. We have distinguished between validated and unvalidated SSNs; we have required that gender and age match across the survey data and the administrative data; we have dropped all imputed data; and we have dropped all movers.

Furthermore, we have used the administrative data to exactly mimic the CPS concept of program participation—Medi-Cal or welfare in any month in the previous calendar year (January to December). While the CPS data refer to the previous year, aggregate Medicaid and welfare tabulations from administrative data are usually published as monthly totals. With movement onto and off Medicaid/welfare, there is no direct relation between the counts for the individual months and CPS concept—the total number of individuals enrolled in the program at any point in a given calendar year. This difference in concepts of participation is not an issue in our analysis. We have the monthly MEDS data and can tabulate it so that it exactly corresponds to the CPS question (i.e., any Medi-Cal in the previous calendar year and any welfare in the previous calendar year).

From these ideal data, we want to compute “behavioral” mis-reporting rates (in contrast to the “imputational rates” that we define in the next chapter). Formally, the mis-reporting rates are:

$$\rho_{FP}^b = \frac{FP}{TN + FP} \qquad \rho_{FN}^b = \frac{FN}{TP + FN}$$

On the left we have the behavioral false positive rate. The denominator is all people who are not truly enrolled, the sum of true negatives and false positives. The numerator is the number of people who are not enrolled, but report that they are. Thus, the ratio is the fraction of people who are not enrolled, who report falsely that they are. On the right is the behavioral false negative rate. By analogy, it is the fraction of people who are enrolled (i.e., the sum of true positives and false negatives), who falsely report that they are not enrolled. The observed net under-reporting suggests that the false negatives are the more common group.

Table 3.3 reports these behavioral rates. The off-diagonals give the false reporting rates; i.e., the percentage of people truly in a given status, who give each possible response. Despite our efforts to develop a sample that should maximize the congruence between the survey and administrative data, the congruence of the two reports is distressingly poor, as can be seen in Table 3.3. We perform the analysis separately for each of the programs of interest: Any Medi-Cal and welfare. There are nine possible outcomes for each individual (three possible MEDS outcomes x 3 possible CPS outcomes, which are summarized in Table 3.3 below.

Table 3.3.
Congruence in the Matched Sample
between MEDS and CPS Data

		CPS			
		N	MO	W	T
MEDS	N	97.7%	1.8%	0.5%	89.8%
	MO	36.8%	59.3%	3.9%	4.9%
	W	20.4%	32.0%	47.6%	5.3%
	T	90.6%	6.2%	3.2%	100.0%

Notes: Entries: row percents (i.e., rows sum to 100%), except for last row and column ("T" for total) that give overall percentages. Computed from the "matched sample" (i.e., adults who provide an SSN).

N-No Medi-Cal, MO-Medi-Cal Only, W-Welfare, T-Total.

Then, the first column of Table 3.3 suggests that few people who do not have Medi-Cal report that they do have Medi-Cal (2.3 percent = 100.0 percent - 97.7 percent). The same is not true for those with Medi-Cal. Less than two-thirds (59.3 percent) of those with Medi-Cal Only actually report having Medi-Cal Only in the CPS (i.e., Medi-Cal, but not welfare); and less than half (47.6 percent) of those with welfare actually report having welfare in the CPS. These are distressingly low rates of congruence, particularly in a sample designed to only include the highest quality matches.

The overall congruence of reports of Medi-Cal is higher. Nearly three-quarters of those with Medi-Cal report that they have Medi-Cal (computed from the second and third rows of Table 3.2 by converting to the unconditional percentages, combining the Medi-Cal Only and Welfare groups and computed error rates). This divergence appears to be explained by those with welfare. While more than half of them report not having welfare (52.4 percent), most of them do report having Medi-Cal. Thus, the overall reporting of Medi-Cal is not as bad as considering welfare and Medi-Cal-Only separately would suggest. This also appears to explain why the raw reporting rates for Medi-Cal Only (i.e., the ratio of the CPS counts to the MEDS counts; see Table 2.3) are so much higher than the reporting rates for Medi-Cal or

Welfare. It is not that reporting rates are higher for Medi-Cal Only; it is that many of those reporting Medi-Cal only in the CPS actually have welfare as well.

The implications of Table 3.2 for reporting rates are relatively subtle. The previous paragraph considered reports conditional on the truth (as recorded in the MEDS). Some of the errors are offsetting. If instead, we compare the gross rates (i.e., the ratio of the CPS total to the MEDS total), the reporting rates are higher. The gross reporting rate for Medi-Cal is 92 percent $((6.2\%+3.2\%)/(4.9\%+5.3\%))$. The gross reporting rate for welfare is much lower at 60 percent $(3.2\%/5.3\%)$. For Medi-Cal Only, the gross reporting rate implies over-reporting of Medi-Cal Only at 128 percent $(6.2\%/4.9\%)$.

These gross reporting rates in the validated years are higher than for the unmatched comparisons that implicitly include both the matched and unmatched samples (about 70 percent of Medi-Cal, about 45 percent for welfare, and about 85 percent for Medi-Cal Only). Thus, while the matched data capture some of the under-reporting, it seems likely that there are additional considerations leading to under-reporting in the unmatched sample. The simplest explanation is higher false negative rates; people who do not provide an SSN are more likely to not report enrollment, even when they are enrolled. Below, we implement an algorithm consistent with that simple explanation.

Other explanations would consider not differential reporting, but differential coverage of the CPS. Even if the CPS survey process was perfect, the CPS only attempts to interview those in the non-institutional population. Anyone receiving Medi-Cal and in an institution would be in the MEDS count, but not in the CPS count. To address this concern, we delete everyone over 65 from both counts. This should eliminate most of the institutionalized population. The size of the remaining institutionalized population is unclear. It seems likely that our inability to better exclude the institutionalized population explains a considerable portion of the apparent excess under-reporting of Medi-Cal among those who do not provide an SSN. We return to this issue in Chapter 4. This issue is unlikely to affect our analysis of the matched sample or of welfare.

Similarly, while the CPS sampling frame is the non-institutionalized population, as with any survey, some people are missed. The CPS adjusts for such failure to interview using control totals derived from the Census that stratify on region, gender, and age. It seems likely, that within these cells, those with Medi-Cal are less likely to be interviewed. If this supposition is correct, then the CPS, even with adjustments, will under-report Medi-Cal enrollment. Again, the methods we propose below will correct for this weighting error, at least at the aggregate level.

Who Mis-reports?

Mis-reporting is not random in the population. Table 3.4 reports eight mis-reporting rates. Behavioral false negatives (the probability of answering “not enrolled” in the CPS, given that the MEDS indicates enrollment); behavioral false positives (the probability of answering “enrolled” in the CPS, given that the MEDS indicated not enrolled); imputational false negatives (the probability of being enrolled according to the MEDS, given answering “not enrolled” in the CPS); and imputational false positive (the probability of not being enrolled in the MEDS, given answering “enrolled” in the CPS). (See below for a discussion of the distinction between behavioral and imputational regressions.) To gain sample size, these tabulations pool the validated and unvalidated years.

The first column gives stratified behavioral false negative rates for Medi-Cal. Overall, about 28 percent of those with Medi-Cal report in the survey that they do not have Medi-Cal. In general, false

negative rates are lower for those with higher Medi-Cal coverage rates (e.g., well below the poverty line, below the poverty line, high school drop outs, single women with children), and higher for those with lower Medi-Cal coverage rates (e.g., some college).

Finally, note the very high negative rates for those with other health insurance. Apparently people who have both other health insurance (as reported in the CPS) and Medi-Cal at some time during the year are much less likely to report their Medi-Cal (as reporting in the MEDS).

These patterns would be consistent with a stigma story. People who are (are not) “expected” to have Medi-Cal would be less (more) reluctant to so report. It would also be consistent with an intensity story. As we will see below, people who are enrolled in Medi-Cal for more of the year are more likely to report Medi-Cal in the survey. People with higher Medi-Cal coverage rates are also more likely to be enrolled for more of the year or at the time of the interview.

Turning to the year patterns, the false negative rate appears to have dropped sharply (from 30 percent to 22 percent) with the survey changes of 1995. However, from 1997 the false negative rate was again over 30 percent.

Table 3.4.
Mis-Reporting Rates (in %)

Sample	Behavioral				Inputational			
	Medi-Cal		Welfare		Medi-Cal		Welfare	
	F-Neg	F-Pos	F-Neg	F-Pos	F-Neg	F-Pos	F-Neg	F-Pos
Overall	28	2	53	1	3	22	3	21
Male	29	2	76	0	2	27	2	22
Hispanic	33	5	59	2	6	28	5	31
Black	29	5	42	2	7	23	5	21
HS Drop Out	26	6	60	2	8	22	8	25
Some College	33	1	42	0	2	25	1	17
FPL<0.5	26	13	63	5	13	26	16	32
0.5<FPL<1.0	19	13	47	5	13	19	16	20
1.0<FPP<1.5	21	10	48	3	11	19	12	19
1.5<FPP<2.0	24	8	50	3	10	19	9	20
Kids in Household	29	3	51	1	5	21	5	20
Single Female w/Kids	22	8	39	5	11	16	12	20
Other Health Insur.	43	1	72	0	2	31	1	31
Welfare	---	---	---	---	---	15	3	21
Medi-Cal	---	---	40	13	---	22	28	21
1990	25	2	46	1	2	27	2	21
1991	30	3	49	1	3	28	3	20
1992	24	3	47	1	2	26	2	25
1993	26	3	49	1	3	24	3	25
1994	30	3	53	1	4	26	3	19
1995	22	2	46	1	3	21	3	23
1996	21	2	48	1	3	18	3	19
1997	30	3	52	1	4	22	3	18
1998	31	2	59	1	4	20	3	19
1999	33	1	59	0	4	14	3	12
2000	33	2	71	1	4	18	3	28

Source: Weighted tabulations from RAND PS-MEDS Match file.

Note “F-Neg” – False Negative, “F-Pos” – False Postive

See text for definition of Behavioral/Reporting False Negative/Positive

Medi-Cal behavioral false positive rates are much lower than the corresponding false negative rates, 2 percent vs. 28 percent. The relative rates are the converse of those for false negatives. Those who are more likely to be on Medi-Cal (minorities, high school drop-outs, in poverty, single women with children) are more likely to report that they have Medi-Cal when they do not. Those who are less likely to be on Medi-Cal (some college, other health insurance) are less likely to report that they have Medi-Cal when they do not. In the next section, we report results consistent with errors in the timing of reporting of coverage.

The reports for welfare are similar. False negative rates are clearly higher for welfare than for Medi-Cal (53 percent vs. 28 percent). People are less likely to report welfare when they have it. The year results suggest a sharply increasing trend at the end of the period covered. This is the period of welfare reform. In general, welfare reform policies were designed to make welfare receipt less socially acceptable. Consequently, a stigma story seems plausible.

It is also possible that these trends were caused by respondent confusion. The names of welfare programs changed. Most states now have a new name for their welfare program (e.g., "CalWORKs" in California and even county-specific names). In some states, contact with the welfare program has shifted to an employment office or to a private sector provider. This later factor, however, is unlikely to be a major issue in California. While some welfare-to-work services in California are provided by other government agencies (e.g., school districts, community colleges, or Workforce Investment Act agencies) or private sector firms, eligibility operations (including monthly or quarterly status reports and annual in-person redeterminations) continue to be handled by county welfare departments in their offices.

As in the Medi-Cal case and also consistent with a stigma story, some college and other health insurance makes not reporting welfare more likely; single woman with children and poverty makes not reporting less likely. The anomaly is deep poverty that makes not reporting much more likely. The explanation appears to be that the welfare question is part of the sources of income battery and welfare is considered income for the purposes of computing the poverty rate. Thus, if one is on welfare and reports the welfare payment, it is difficult to be in deep poverty (income below half the corresponding poverty line). Again, the false positive rates for welfare are the mirror image of the false negative rates for welfare. Those with a higher (lower) probability of receiving welfare are more (less) likely to report receiving welfare, even when they do not.

"Behavioral Regressions" and "Imputational Regressions"

Clearly from both a statistical perspective and a disclosure perspective, the sample sizes are only barely large enough to support contingency table analyses. Furthermore, we would like to understand how reporting has shifted over time and the potential effects of the changes in the CPS questions, most notably those in 1996 and 2000. Appendix B reports a total of 8 regressions, corresponding to the eight columns of Table 3.4. Appendix B also provides details of the variable specifications, the stepwise strategy, and the regression coefficients.

These results are crucial for the imputation models that follow, but they are otherwise difficult to interpret. Crucially for our purposes, they confirm the aggregate data and show no trend in the imputation false negatives for Medi-Cal but confirm a strong increase in false negatives for welfare. This is consistent with the tabulations from the unmatched data in Table 2.2.

Understanding the Reporting Errors

Table 3.3 tabulates the MEDS data according to whether there was any program enrollment in the previous year. Standard cognitive approaches to survey response (Sudman and Bradburn, 1973; Groves, 1989) would suggest that a positive Medi-Cal response is more likely the more Medi-Cal/welfare receipt there was and also if there is Medi-Cal receipt in the survey month.

Our detailed MEDS data allow us to compute months of enrollment last year and enrollment in March of this year. The observation counts at each month are too small to allow reporting the raw rates. Instead for Medi-Cal, Figure 3.1 reports: (1) the raw rates for March of this year, and all of last year; and (2) not March of this year and none of last year. For the other combinations, we report the results of logit polynomial regressions.¹⁶ Reporting the regression results also smoothes out some of the sampling variation.

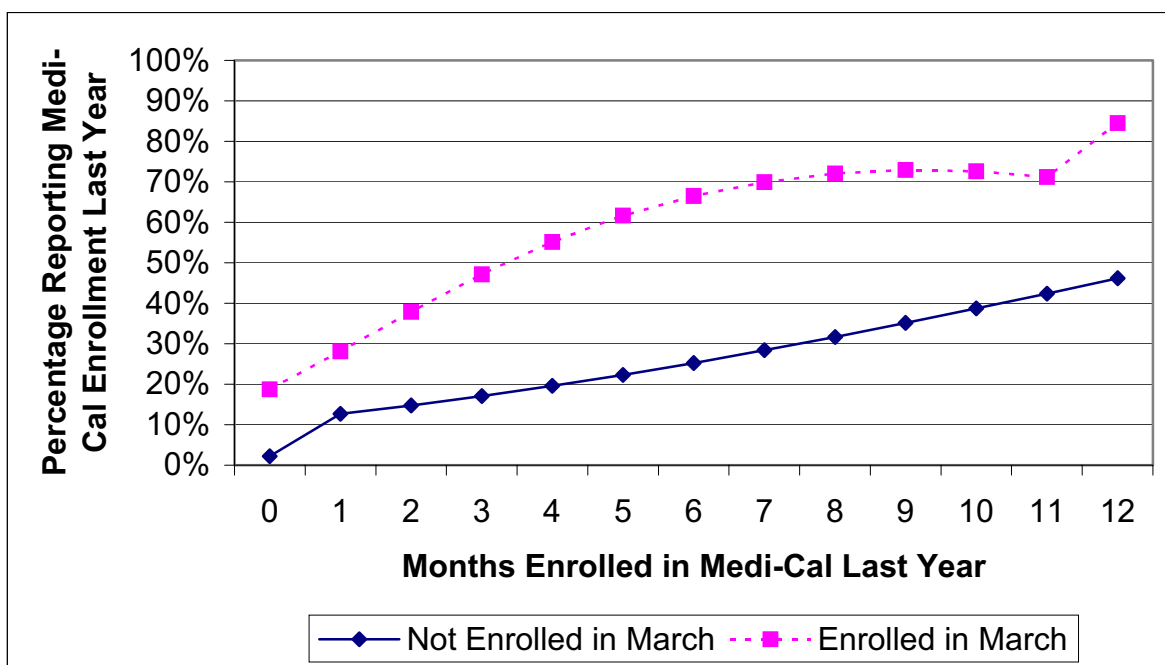


Figure 3.1—CPS Reporting of Medi-Cal Given MEDS Pattern of Receipt

If reporting were perfect, the points for “0” would be zero and the other points would be 100 percent. Instead, we observe a clear dose-response relationship. The more months of enrollment last year recorded in the MEDS, the more likely a person is to report program enrollment in the CPS. People on Medi-Cal all of last year and in March of this year, report Medi-Cal at about 85 percent. People on in March, but not on all of the previous year, are less likely to report Medi-Cal, with reporting rates varying from 30 to 70 percent.

¹⁶ The probit analysis takes as its dependent variable the percentage of people reporting enrollment from the MEDS, given their actual status in March of the survey year (i.e., the vertical axis of Figure 3.1) and as its independent variables polynomials in months of actual enrollment from the MEDS (i.e., the vertical axis of Figure 3.1). The polynomial is quadratic for those on welfare this March and linear for those not on welfare this March. The analysis uses the CPS sample weights. We then plot the predicted probabilities from that model.

The CPS question specifically asks about enrollment last year, but enrollment in March clearly affects the probability of answering the CPS question about last year positively. The difference ranges from 15 to 30 percentage points. Finally, note that people enrolled in March but not at all last year—who should answer negatively—have a 20 percent probability of answering in the affirmative (i.e., they respond based on their current enrollment status rather than their enrollment status last year). These people appear to explain much of the false positives, people who report enrollment in the CPS but who are not actually enrolled (in the past year) according to the MEDS.

Figure 3.2 reports the same tabulations and logitlogit regression predictions for welfare. The patterns are similar. People enrolled in March are more likely to respond positively. Even some (about 5 percent) of those enrolled in March, but not last year, respond positively, comprising many of the false positives. Both for those enrolled in March and for those not enrolled in March, a positive response is more likely the more months of enrollment in the past year. However, even for people enrolled in March and all of last year, only slightly more than half report their welfare participation in the CPS.

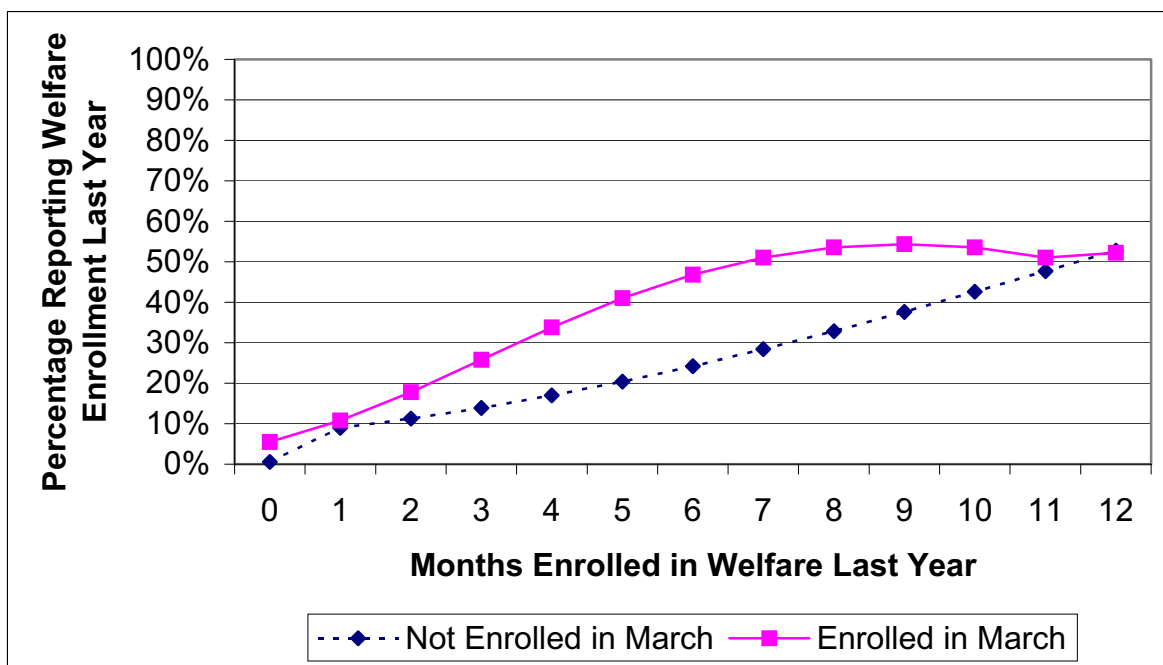


Figure 3.2 – CPS Reporting of Welfare Given MEDS Pattern of Receipt

The CPS also includes a question about months of receipt of Medicaid and welfare last year; however, the individual cells are too sparse to allow reporting. We do find that there is some correlation between the MEDS and CPS data, but the correlation is only moderate (0.58 for Medi-Cal, 0.57 for welfare) and the mean numbers of months reported are very different. Consequently, we conclude that the monthly data in the CPS are not very useful.

Some of the other welfare false positives are also understandable. The CPS distinguishes between AFDC/TANF/CalWORKs (a state program) and SSI (a federal program). Both programs provide cash and both automatically confer Medi-Cal eligibility. Thus, it would not be surprising if some people on SSI report to the CPS that they are on welfare. We would code those people as “false positives.” Because SSI confers Medi-Cal eligibility, we can identify SSI in MEDS and, therefore, in the matched data.

Tabulations on the validated years suggest that a sizable fraction of the welfare false positives in fact have SSI. Cell counts are too small to support more precise estimates or their release. Even with these explanations, some of the welfare false positives are still unexplained. Based on these results for SSI, however, it seems likely that some of the remaining false positives are related to receipt of General Relief (which we can not identify in our data).

Unfortunately, we cannot use similar methods to better understand the false positive reports for Medi-Cal, because we do not have similar administrative data that identify other sources of public insurance. By analogy to the results for welfare, however, it seems likely that some other public health insurance programs are being reported as Medicaid, again inducing false positive responses. For example, in California, Healthy Families is administered outside the Medicaid system and not recorded in MEDS. Thus, some of the Medi-Cal Only false positives are probably Healthy Families enrollees. However, given that we have dropped everyone under age 15 and Healthy Families in California does not cover adults, this is unlikely to be a major factor. More recently, several counties have put in place public, county-level Medicaid-like programs. Although these programs were implemented after the period our data cover and, thus, cannot explain the false positives in our data, this may be an issue in CPS interviews in later years, where such programs may incorrectly be reported as Medicaid.

Another possible explanation for the false positives in both of the programs is migration between states. As we discuss in detail at the end of Appendix A, we control for migration, but those controls are incomplete. For example, some people who we have classified as California residents in fact spent part of the reference year in another state. If they were enrolled in Medicaid/welfare there, we would have incorrectly labeled them as false positives.

We note, in contrast, that the false negative rates (i.e., the percent of people enrolled in the program based on the MEDS that do not report participation in the CPS) seem much too high to be explained away by any of these factors.

CPS Reference Period

Figure 3.1 and Figure 3.2 suggest that enrollment in the following March increases the probability of a positive response. Of particular note are the people who were not enrolled at any point during the last year, but are enrolled in March of this year. They should give a negative response to a question about enrollment last year. For Medi-Cal, about 20 percent of this group gives a positive response; for welfare the comparable figure is about 5 percent. These percentages are about half of the percentages for people enrolled in March of this year, but only enrolled one month of the previous year. These results are consistent with Swartz (1986) and Sudman, Bradburn, and Schwarz (1996) who argue that the CPS responses should be interpreted as referring not to the previous calendar year, but instead to the interview month. CBO seems to accept this argument (Bilheimer, 1997; CBO, 2003). Lewis, Ellwood, and Czajka's (1998, p. 27) in their review conclude:

The CPS is designed to measure the number of individuals uninsured throughout a given year. Yet most researchers believe the CPS estimates of the uninsured represent a mix of those uninsured *throughout* [emphasis in the original] the previous year and those uninsured at a point in time [i.e., as of the interview].

Interpreting the CPS responses as referring, not to the previous year, but instead to the interview month would in general yield both lower Medi-Cal enrollment rates (fewer people are enrolled in a given month than in an entire year) and higher uninsurance rates (fewer people are uninsured in an entire year than in a given month).¹⁷

From our matched sample, we can explore this question directly (Table 3.5). We divide our sample into four groups, by their true (i.e., MEDS) enrollment status in the survey month (yes/no) crossed with their true enrollment status in the previous year (yes/no). As expected, those enrolled both last year and in the survey month are most likely to report in the CPS that they are enrolled (79 percent for Medi-Cal, 50 percent for welfare), while those enrolled at neither time are least likely to report in the CPS that they are enrolled (2 percent for Medi-Cal, 1 percent for welfare).

Table 3.5.
CPS Reference Period

MEDS		Medi-Cal			Welfare		
Last Year	Survey Month	CPS='Y'	CPS='N'	% Y	CPS='Y'	CPS='N'	% Y
Y	Y	6.85%	1.80%	79%	1.91%	1.91%	50%
Y	N	0.58%	1.56%	27%	0.46%	1.16%	28%
N	Y	0.08%	0.43%	16%	0.00%	0.19%	0%
N	N	1.99%	86.71%	2%	0.56%	93.81%	1%
Overall		10%			3%		

Notes: First two columns give MEDS status ("Last Year" – any enrollment in the previous calendar year; "Survey Month" – enrollment in the March of the current year). Columns labeled "CPS-'Y'" and "CPS-'N'" give the unconditional probability of enrollment in Medi-Cal/Welfare. Column labeled "% Y" gives the probability of a "Y" response in the CPS given the MEDS data for the survey month and the previous calendar year (i.e., within the row).

Our interest focuses on the cases where the status last year diverges from the status at the interview. The CPS question explicitly asks about the previous year, so the second row "should be" 100 percent and the third row "should be" zero. In fact, the fraction responding in the affirmative is higher for the second row than for the third row (for Medi-Cal, 27 percent vs. 16 percent; for welfare 28 percent vs. 0 percent). Thus, CPS respondents are more likely to interpret the question as referring to last year, than as referring to the interview month.

Table 3.4 has several other implications. First, the distinction between current month and previous year cannot be very important. Very few people have different statuses according to the two concepts. Second, once the two concepts diverge, neither response is of high quality. Unless a respondent was enrolled both last year and in the survey month, a "No" answer is the most likely. Finally, enrollment last year but not in the survey month is much more common than enrollment in the survey month but

¹⁷ See Bennefield (1996b) who argues that the problem is generic under-reporting because of the length of the recall period. See also Fronstin, 1996, on dual coverage; Kronick, 1991, on private health insurance; Beauregard et al., 1997, comparing to MEPS results; Bennefield, 1996c, using CPS experimental questions; and Long and Marquis, 1996, comparing to RWJF survey.

not last year. This is as would be expected in a time stationary environment (the probability of receiving in at least one of twelve months is higher than the probability of receiving in any give month).

Time Trends

The previous analyses have pooled results across the (validated) years for which we have matched data. However, two factors suggest the importance of studying change over time. First, based on the aggregate data there is some evidence of increases over time in under-reporting. Second, the Census Bureau has altered the CPS questions, most notably in 1996 and 2000, which appears to have increased reports of health insurance coverage.

Unfortunately, the matched samples in each year are small, about 5,000 adults per year. Medi-Cal enrollment rates average about 10 percent, about 500 people in the CPS; welfare and Medi-Cal only rates are about half that. Furthermore, false positives are rare, often a few percent or less. As a result, we cannot provide descriptive evidence on time trends in the probability of false positives and false negatives.

CPS Imputations

The previous analyses in this chapter have used what we called the matched sample. To ensure that any lack of congruence resulted from true response errors, we deleted all imputed data from the analysis.

However, we can and did perform the match for all people for whom an SSN was available. We will use this sample to analyze the quality of CPS imputation of program participation. As in the main analysis, we drop those without SSNs, apparently bad matches, and movers. Table 3.6 reports the congruence of the MEDS data and CPS imputed responses. The table has three panels. The first panel considers what the Census refers to as “allocation” of welfare (i.e., imputation for item non-response). This allocation is done by the standard Census hot-deck procedure (i.e., missing data are imputed based on randomly choosing a case with non-missing data on the variable of interest and similar demographic and economic values on the non-missing data). Similarly, the second panel considers the “allocation” of Medicaid responses for item non-response. Finally, the third panel considers logical imputation of Medicaid. The CPS automatically imputes Medicaid to anyone who reports the receipt of welfare or SSI.

Table 3.6.
MEDS Data for CPS Imputed Records

CPS Value	MEDS Data		Total
	No	Yes	
Welfare Allocated			
No	94.0%	2.8%	96.8%
Yes	0.76%	2.5%	3.28%
	94.6%	5.4%	100.0%
Medi-Cal Allocated			
No	56.0%	30.3%	86.3%
Yes	6.2%	7.5%	13.7%
	62.2%	37.8%	100.0%
Medi-Cal Logically Imputed			
Yes	51.2%	48.8%	100.0%
Total	51.2%	48.8%	100.0%

Note: Cells are percent within the panel. Verified years only.

Because the CPS does not impute for program participation for many people, the sample sizes for the estimates presented in Table 3.5 are small (about 100 cases are imputed over all of our data). Consequently, these results need to be treated with caution. With that caveat, the results suggest that the hot-deck algorithm is under-estimating program enrollment. For welfare, the difference is small. The MEDS suggests that about 5.4 percent of the allocated cases are enrolled; the Census imputes welfare to only 3.2 percent of them. For Medi-Cal the differences are larger. In MEDS, 37.8 percent of the cases are enrolled in Medi-Cal whereas the CPS imputes only 13.7 percent.

The logical imputations add Medicaid to children whose parents report Medicaid. For older children, we can check this imputation against the MEDS. We find that only 38.8 percent (partially imputed to Medi-Cal) actually have Medi-Cal. The imputation is making things worse.

In the matched sample, the logical imputations are about two percent of all cases and the allocations about one percent of all cases. Relative to perfect imputations, the incorrect logical imputations therefore raise the Medi-Cal enrollment rate by about one percentage point. Relative to perfect imputations, the incorrect allocations lower the Medi-Cal enrollment rate by about a fifth of a percentage point.

For welfare, they represent less than a fifth of a percentage point. The effect on overall enrollments of any imputation errors is therefore trivial. In fact, the imputations of welfare participation appear to be quite good.

These results suggest that people who do not answer the Medicaid questions are substantially more likely to have Medicaid/Medi-Cal than the demographically similar households the CPS hot-deck procedure is using for its imputations. In short, item non-response is not random, but the effect on total estimated enrollment is trivial. In contrast, the Medicaid logical imputations are wrong about half the time, increasing estimated Medi-Cal enrollment by about one percentage point, which is about ten percent of true Medi-Cal enrollment. Some additional attention to the Medicaid logical imputations may be appropriate.

Conclusion

This chapter has considered congruence between the CPS and MEDS data in the best possible matching sample. Even in this sample, the level of congruence is distressingly low and appears to be getting worse over time for welfare.

4. Extrapolating to the Full Data

The previous chapter analyzed the congruence of responses among those who provided an SSN (what we referred to as the matched sample) and met a set of sample inclusion criteria. However, as Table 3.1 notes, many people do not provide an SSN. Furthermore, not providing an SSN is differential. People who are more disadvantaged are less likely to provide an SSN, but people who are on welfare are slightly more likely to provide an SSN.

The fundamental problem that is addressed in this chapter is that we do not have SSNs for about half of the sample. We do not want to assume that the responses in the unmatched sample are perfect. Instead, we want to use information from the reporting errors in the matched sample (where we have the MEDS information, treated as truth) to perform better imputations of program enrollment in the unmatched sample. The basic idea is that individuals with characteristics associated with under-reporting in the matched sample are more likely to under-report in the unmatched sample. We estimate a logistic regression model of such reporting errors (both under-reporting and over-reporting) on the matched sample. We then use that model to multiply impute a true response in the unmatched sample; where by multiple imputation we mean that we assign a probability of each response to each individual based on the regression model.

In practice, we have one more piece of information. We can estimate the total number of people in the unmatched sample who are enrolled in a program. To do so, we take the total estimates from the administrative data and subtract the estimates of enrollment in the matched sample (i.e., we use the CPS weights and the MEDS/administrative data information). Our logistic regression models in general under-predict the number of program enrollees in the matched sample. We therefore append a multiplicative adjustment factor. The effect of that adjustment factor is to force the imputed number of program enrollees to exactly match the administrative totals.

The balance of this chapter provides a precise mathematical discussion of the problem and our approach. The discussion in this chapter is formal and technical. Many readers will want to skip to the next chapter where we provide the substantive results.

The Identification Problem

We can conceptualize the CPS matching problem as a table including eight “cells,” in terms of total weighted counts. The columns distinguish whether the individual is on the program according to the MEDS (i.e., YES/NO). The rows distinguish both the CPS response and whether the record has an SSN (so it is potentially matchable). The letters name the cells to ease the discussion below.

CPS	SSN		MEDS		Total
			YES	NO	
	Present	YES	$A=TP_S$	$B=FP_S$	$C=Y_S$
		NO	$D=FN_S$	$E=TN_S$	$F=N_S$
	Absent	YES	$G=TP_A$	$H=FP_A$	$I=Y_A$
		NO	$J=FN_A$	$K=TN_A$	$L=N_A$
		Total	$M=Y_M$	$N=N_M$	$O=T$

Thus, the subscripts are:

- “S” – SSN present (i.e., a match was in principle possible; in practice, we drop the bad matches as well);
- “A” – SSN absent (i.e., a match is not possible);
- “M” – MEDS.

And the other codes are:

- TP – true positive;
- TN – true negative;
- FN – false negative;
- FP – false positive.

And finally:

- Y – “Yes” (on Medi-Cal/welfare);
- N – “No” (not on Medi-Cal/welfare);
- T – Total.

We treat the MEDS data as “truth.” Thus, our goal is to use the MEDS data to “fix” the CPS data. From the records that provided SSNs, we know TP_S , FP_S , FN_S , and TN_S . So, we simply adjust the CPS answers to align with the MEDS answers.

The challenge, therefore, is the unmatchable data—those records for which no SSN was available in the CPS.¹⁸ For those records, we only know the row totals— Y_A , N_A —and the column totals by subtraction from the “S” sample— $TP_A+FN_A=Y_M-TP_S-FN_S$ and $FP_A+TN_A=Y_M-FP_S-TN_S$. However, there is some additional—we will see, not quite enough—information from the matched sample.

To understand our approach, begin by formally defining the *imputational false positive rate* and the *imputational false negative rate* as:

¹⁸ Any errors in coverage of the CPS are likely to have their effects at this stage. Consider the possibility that the CPS systematically and disproportionately misses (i.e., fails to interview) those with welfare/Medicaid and those errors are not corrected by CPS reweighting (for region, urban/rural, gender, age, race-ethnicity). The MEDS and the MEDS counts would include these people. The CPS and the CPS counts would not. Our approach will attribute this to under-reporting. In fact, the problem would be, not that respondents responded falsely, but that people with welfare/Medicaid were not interviewed.

This problem is likely to be most severe among the institutionalized population. They are not included in the CPS sampling frame. They are in our MEDS counts. We partially address this issue by dropping those over 65. It seems likely that a better estimate of the number of people with Medi-Cal in institutions would considerably lower our estimated false negative rates in the unmatched population relative to the false negative rates in the matched population. This is the α we define below. This issue probably explains much of why α is greater than 1.

$$(4.1) \quad \rho_{FP}^i = \frac{FP_S}{TP_S + FP_S} \qquad \rho_{FN}^i = \frac{FN_S}{TN_S + FN_S}$$

where the “i” is for imputation and these are the rates with respect to the CPS answers (as opposed to the behavioral rates in terms of the true behavior as measured in the MEDS that we also considered in the previous chapter).

These imputational rates are in contrast to the behavioral rates of the previous chapter. The tabulations there addressed the behavioral question. Given the true status, what is the probability of a false response? This is not a useful concept for imputation. In the CPS, we observe the potential false response and want to infer the true status. To do so, we want the imputational rates, i.e., the probability that the true status is different than the observed response, given the observed response. The two sets of rates are exactly related. From a complete 2x2 contingency table (i.e., TP, FP, TN, FN), we can compute both sets of rates. From one set of marginals and one set of rates, we can recover the other set of rates. Which rate is more insightful depends on whether we are addressing behavioral questions (as in the last chapter) or imputational questions (as in this chapter).

Then, if we knew these imputational error rates, we could probabilistically impute the data. We would create two pseudo-observations for each observation (dividing the sample weight between the pseudo-observations). So, for example, if an observation reported “Y” in the CPS, that observation would be assigned a “Y” with probability $1 - \rho_{FP}^i$ and “N” with probability ρ_{FP}^i . Similarly, if an observation reported “N” in the CPS, that observation would be assigned a “N” with probability $1 - \rho_{FN}^i$ and a “Y” with probability ρ_{FN}^i .

We do not know these rates in the unmatched sample. Furthermore, the rates from the matched data are not directly applicable in the unmatched data. If the rates from the matched sample applied in the unmatched sample, then applying those rates to the unmatched data would recover the actual number of people on Medi-Cal/welfare in the MEDS, i.e.:

$$(4.2) \quad Y_M = TP_S + FN_S + Y_A (1 - \rho_{FP}^i) + N_A$$

However, we have already noted that the under-reporting in the matched sample is not large enough to explain the under-reporting in the full sample.

We have a fundamental non-identification problem: One equation for Y_M and two unknowns—the rates in the unmatched data. Setting one of the rates fixes the other rate.

Given that in net we have under-reporting of Medi-Cal/welfare and false positives are rare (and relatively stable through time), we adopt the simplest rule. We use the false positive rate from the matched data in the unmatched data. We then adjust the false negative rate (by a multiplicative factor, “ α ”) until the implied total count of people on Medi-Cal/welfare in the CPS equals the count in the MEDS (assumed to be truth).

$$(4.3) \quad Y_M = TP_S + FN_S + Y_A (1 - \rho_{FP}^i) + N_A \alpha$$

The left-hand side is the “true” number of individuals on Medi-Cal/welfare from the MEDS. The right-hand side is the “fixed” number of individuals on Medi-Cal/welfare in the CPS. Considering each of those terms in turn:

- $TP_S + FN_S$: The number of people who have Medi-Cal/welfare in the matched sample (true positives plus false negatives).
- $Y_A(1 - \rho_{FP}^i)$: The number of people who report having Medi-Cal/welfare who actually do. We know the number of people who report having Medi-Cal/welfare in the unmatched sample. We estimate the number of these people who actually do have Medi-Cal/welfare using the imputational false positive rate from the matched sample. This is the identifying assumption.
- $N_A\alpha$: The number of people who report not having Medi-Cal/welfare who actually do. Again, we know the number of people who report not having Medi-Cal/welfare in the unmatched sample. Finally, α gives the probability of a false negative in the unmatched data.

Solving for α , the false negative rate in the unmatched sample yields:

$$(4.4) \quad \alpha = \frac{Y_M - TP_S - FN_S - Y_A(1 - \rho_{FP}^i)}{N_A}$$

Except for the false positive rate, each of the terms on the right side is observable. In the numerator, the first term is the number of people on Medi-Cal/welfare from the MEDS. The second and third terms are the number of people on Medi-Cal/welfare from the MEDS in the matched sample. The fourth term is the product of the number of number of people in the unmatched sample who claim to have Medi-Cal/welfare and the fraction of them who are estimated to actually have Medi-Cal/welfare (i.e., one minus the imputational false positive rate). The denominator is the (weighted) number of people in the CPS sample who do not provide an SSN who claim not to have Medi-Cal/welfare. The false positive rates for the unmatched data are not observed, but by assumption we use the value estimated from the CPS. Since the CPS is a sample, each of these concepts should be weighted.¹⁹

Imputing the Data

This analysis of identification suggests that we are missing one piece of information. However, once we assume that the false positive rate is common in the matched sample and the unmatched sample, we can solve for α . Then, knowing α is enough to solve for each of the cells:

¹⁹ We note that this is the analysis considering the concepts (Medi-Cal, welfare) separately. It would also be of interest to impute jointly welfare and Medi-Cal. Table 3.3 and the probit regressions reported in Table 3.8 provide the inputs for such an analysis.

We do not perform the full imputation here. The actual imputation would be more complicated than the single imputation attempted here. The single imputation considered here is for a 2x2 table, with two error rates. Fixing one of the error rates is enough to allow computation of the other one from the data.

In contrast, the joint response problem is a 3x3 table, with six distinct error rates. We need to fix four of them to be able to compute adjustment factors for the last two. By analogy, with the approach in the body of the paper, it would be natural to assume that the three upcoding error rates are common. However, that is not sufficient. We still need to fix either the $P[W|N]$ or $P[W|MO]$. We have seen that both of these errors are common and changing over time, so it is not clear how to proceed.

$$\begin{aligned}
 (4.5) \quad TP_A &= Y_A(1 - \rho_{FP}^i) \\
 FP_A &= Y_A \rho_{FP}^i \\
 FN_A &= N_A \alpha \\
 TN_A &= N_A(1 - \alpha)
 \end{aligned}$$

Cell counts for the terms in individual years are often too small to allow public release. However, the totals over the whole 11-year period are releasable. To understand our methods, Equation 4.6 and Equation 4.7 show the actual numbers for Medi-cal and welfare respectively, summing over all 11 years (rounded to hundreds of thousands).

$$\begin{aligned}
 (4.6) \quad \alpha &= \frac{Y_M - TP_S - FN_S - Y_A(1 - \rho_{FP}^i)}{N_A} \\
 &= \frac{31.4 - 10.6 - 5.0 - 6.1}{4.0} = 2.4 \approx 2.4
 \end{aligned}$$

$$\begin{aligned}
 (4.7) \quad \alpha &= \frac{Y_M - TP_S - FN_S - Y_A(1 - \rho_{FP}^i)}{N_A} \\
 &= \frac{13.9 - 3.6 - 4.5 - 1.5}{3.1} = 1.4 \approx 1.4
 \end{aligned}$$

Thus, over the full 11 years, the MEDS has 31.4 million adults on Medi-Cal. The weighted matched CPS data have 10.6 million true positives and 5.0 million false negatives. Using the imputational false positive rates and false negative rates, we would estimate 5.0 million false positives and 6.1 million false negatives in the unmatched sample. To align the CPS totals with the MEDS totals, we need to increase the false negative count by a factor of 2.4 (i.e., from 4.0 to 9.6).

For welfare, the MEDS has 13.9 million adults on welfare. The matched CPS data have 3.6 million true positives and 4.5 million false negatives. Using the imputational false positive rates and false negative rates, we would estimate 1.5 million false positives and 3.1 million false negatives in the unmatched sample. To align the CPS totals with the MEDS totals, we need to increase the false negative count by a factor of 1.4 (i.e., from 3.1 million to 4.3 million).

Stratifying and Adjusting for Covariates

The above analysis is applicable when the population is homogeneous. In reality, the population is heterogeneous. We are able to address this to some extent. We have a small number of variables—calendar year (in principle, also gender and age)—that are measured (nearly) consistently in the MEDS and the CPS. For these variables, we can totally stratify (i.e., we will compute a different value of α for every strata).

In addition to the small number of variables that are common to both data sets, we have many other covariates in the CPS. This allows us to estimate the false negative and false positive rates in the matched

sample, not only in terms of the small set of common variables, but also in terms of the larger number of variables in the CPS alone. We do so using exactly the imputational regressions discussed in Chapter 3.

Using these multivariate models seems particularly important for two reasons. First, many of these variables are likely to be strongly related to Medi-Cal/welfare eligibility and therefore to true Medi-Cal/welfare coverage, e.g., marital status, presence of children in the household, and household earnings. Second, dual coverage (Medicaid and also other, usually private, health insurance) is an issue of substantive interest. As much as possible, we want to correctly impute in the sub-samples with and without private health insurance.

Suppose that within the strata, s , we can assign each individual his/her own ρ (thus the i superscript), then we can write our equation for α as:

$$(4.8) \quad \begin{aligned} Y_{M,s} &= TP_{S,s} + FN_{S,s} + \sum_{j \in Y_A} w_j (1 - \rho_{FP}^i[j]) + \alpha_s \sum_{k \in N_A} w_k \rho_{FN}^i[k] \\ \alpha_s &= \frac{Y_{M,s} - TP_{S,s} - FN_{S,s} - \sum_{j \in Y_A} w_j (1 - \rho_{FP}^i[j])}{\sum_{k \in N_A} w_k \rho_{FN}^i[k]} \end{aligned}$$

In practice, we use the predictions of the imputational logitlogistic regression models from the previous chapter to estimate the ρ s.²⁰

For this project, we have the matched data. However, given the imputational logitlogit model, this approach can also be applied to the CPS public use data by those who do not have the matched data. To see this write

$$(4.9) \quad \begin{aligned} Y_{M,s} &= \sum_{j \in Y_S \cup Y_A} w_j (1 - \rho_{FP}^i[j]) + \gamma_s \sum_{k \in N_S \cup N_A} w_k \rho_{FN}^i[k] \\ \gamma_s &= \frac{Y_{M,s} - \sum_{j \in Y_S \cup Y_A} w_j (1 - \rho_{FP}^i[j])}{\sum_{k \in N_S \cup N_A} w_k \rho_{FN}^i[k]} \end{aligned}$$

Below, we compute α and γ for each strata. Thus, an analyst without access to the matched data could also create an imputed data set.

In what follows, we apply this approach directly to our CPS data. We stratify by year. In practice, the estimates within demographic sub-groups are too small to yield reliable estimates. Table 4.1 presents the resulting estimates for α . For adults, the adjustment factors are quite large in the early years. For Medi-Cal, we need to triple or even quadruple the false negative rates in the early years, suggesting that the unmatched sample is very different from the matched sample. For welfare, despite the fact that the

²⁰ The form of the equation in the text is computationally straight-forward. One could argue that it would be more consistent with the logit modeling strategy to include α inside the logit index. Doing so would require a non-linear optimization to compute α .

under-reporting is absolutely more severe, the unmatched sample is closer to the matched sample. The highest adjustment factors are only slightly greater than two.

Over the 11 years covered by our analysis, the adjustment factors shrink. By 2000, the adjustment factor for Medi-Cal is under 1.5 and for welfare, under 1. It is not clear whether these changes over time result from changes in the CPS instrument or from changes in who is receiving welfare. The large drop in 1995 is consistent with the desired effects of the change in the CPS instrument in that year. The drop in 2000 would also be consistent with the changes in the CPS instrument in that year. Unfortunately, the drop seems to date back to 1999, one year too early.

Table 4.1.
Adjustment Factors α

Year	Adults		Children	
	M	W	M	W
1990	3.0	1.8	1.7	3.3
1991	4.2	2.0	2.3	4.1
1992	3.7	2.2	2.6	4.0
1993	3.2	1.8	3.1	5.2
1994	3.7	1.3	1.0	3.3
1995	2.3	1.3	2.0	3.4
1996	2.2	1.2	1.7	3.0
1997	2.9	1.5	2.2	3.8
1998	2.2	1.3	2.0	3.7
1999	1.4	1.2	1.3	2.8
2000	1.4	0.7	2.4	4.9
Average	2.4	1.4	2.0	3.7

The preceding discussion applies to adults, for whom we potentially have an SSN. We do not have SSNs for any children. Following Census practice, we impute from parents to children.²¹ For Medi-Cal this is consistent with Census's logical imputations. If parents have Medicaid, then children are imputed to have Medicaid. For welfare, this is definitional. Children are never asked about welfare in the CPS. Instead, we impute welfare to both adults and children based on the receipt of public assistance from a welfare program. We then follow the equivalent approach; in other words, we adjust the false negative rate until it aligns the imputed data with the MEDS totals.

We note that the adjustment factors for children are much higher. Furthermore, unlike the adjustment factors for adults, the adjustment factors for children do not fall over time. These adjustment factors are large enough to cast some doubt on the quality of the imputations for children. The adjustments will align the total number of children with the control counts. Our methods impute to

²¹ See for example the March CPS documentation for 1990 (p. 9-8;

<http://www.census.gov/apsd/techdoc/cps/cpsmar00.pdf>): "After data collection and creation of an initial microdata file, further refinements were made to assign Medicaid coverage to children. In this procedure all children under 21 years old in families were assumed to be covered by Medicaid if either the householder or spouse reported being covered by Medicaid (this procedure was required mainly because the Medicaid coverage question was asked only for persons 15 old and over). All adult AFDC recipients and their children, and SSI recipients living in States which legally require Medicaid coverage of all SSI recipients, were also assigned coverage."

children in proportion to the false negative rates. This continues to be a reasonable approach. However, the adjustment factors are so large as to suggest that there is some factor beyond false negatives explaining the under-reporting for children. Whatever it is, the matched data do not identify it.

Conclusion

Given the adjustment factors shown in Table 4.1, we create a multiply-imputed data set. For the matched data, we overwrite the CPS data with the MEDS data. For the unmatched “Yes” responses, we multiply impute based on the false positive rates implied by the logitlogit regression coefficients from the matched sample. For the unmatched “No” responses, we multiply impute based on the product of the false negative rates implied by the logitlogit coefficients from the matched sample and the adjustment factor, α , for this survey year. In practice, we create two data sets, one for the analysis of Medi-Cal (and health insurance) and a second for the analysis of welfare. We do not attempt the full joint imputation of Welfare and Medi-Cal Only.

These multiple-imputation models appear to be the best that can be done with the available data—including the CPS-MEDS match. For those who provide an SSN (our “matched sample”), we use the MEDS data to impute program participation. For those who do not provide an SSN (our “unmatched sample”), we use a multivariate model estimated on those providing an SSN to impute program participation. If the unmatched sample were identical to the matched sample except that they did not provide an SSN, the α s would be close to one. In fact, in most cases, the α s are greater than one, suggesting that people who do not provide an SSN are—even conditional on covariates—more likely to be program participants than those who do not provide an SSN.

In the subsequent chapters we use the multiply-imputed data sets to obtain a better understanding of how the mis-reporting in the CPS can affect different types of analyses. Specifically, in Chapter 5, we examine how mis-reporting of Medi-Cal receipt affects estimates of the number of uninsured in California and in Chapter 6 we look at how mis-reporting of Medi-Cal and welfare participation affect estimates of program take-up.

These substantive estimates in the next two chapters use the results of the model. Some of the α s are so large (well above 2, especially for children) to suggest that the results should be used with some caution. We nevertheless believe that they represent an improvement over estimates that make no correction for under-reporting or even estimates that correct for under-reporting without access to matched data.

5. New Estimates of the Uninsured

Having characterized the under-reporting problem and described our approach to estimating true rates from the matched data, in this chapter we present the first substantive results of this paper – adjusted health insurance rates, overall and by subgroups.

Dual Reporting

To provide improved estimates of the number of people who are uninsured, the crucial issue concerns dual coverage. We know the number of people with Medi-Cal exactly from the MEDS administrative data. However, our matched data provides no new information on who has private health insurance coverage. If no one had both Medi-Cal and private health insurance, we could compute the number of uninsured as the total population less the MEDS estimate of those on welfare and the CPS estimate of those with private health insurance.

However, dual coverage is possible. First, it is possible that a person has both Medi-Cal and other health insurance in a given month. Second, over the course of a year, some months of Medi-Cal and some months of other insurance are even more likely.

On the assumption that CPS responses about other health insurance are correct, we can use our matched sample to tabulate rates of dual coverage.²² In the matched sample, about 23.9 percent of adults who report to the CPS that they have Medi-Cal also report private coverage. (See Table 5.1.) Over the entire sample, the figure for children is also 23.9 percent. For survey year 2000 (the last year for which we have matched data; note that this refers to calendar year 1999), slightly more adults are dually covered (27.6 percent) and slightly fewer children (20.3 percent).

²² This analysis implicitly assumes no reporting error in other health insurance responses. As with all survey data, there is undoubtedly some error in other health insurance responses. Our MEDS data match did not include information allowing us to evaluate such response error rates.

The following considerations may be relevant. First, stigma is unlikely to be a major cause of false negatives. If anything, stigma might induce false positives; i.e., respondents lie by giving the socially more acceptable and empirically more common answer. Second, simple recall errors (e.g., not reporting private health insurance last year, because one does not have it at the interview in March) seems likely.

Table 5.1.
Estimates of Dual Coverage and Uninsurance

	Pooled				2000/1999			
	Adults		Children		Adults		Children	
	DC	UI	DC	UI	DC	UI	DC	UI
Under-reported		4.1%		9.0%		3.8%		9.7%
Method 1		23.5%		17.8%		23.3%		17.8%
Method 2	0.0%	19.5%	0.0%	8.9%	0.0%	19.5%	0.0%	8.1%
Method 3	23.9%	20.4%	16.7%	10.3%	27.6%	20.6%	18.3%	9.8%
Method 4	32.3%	20.8%	34.5%	11.9%	31.2%	20.7%	35.7%	11.5%
Method 5	31.6%	20.8%	34.5%	11.9%	26.4%	20.5%	35.8%	11.5%

Note: DC-estimate of dual coverage used to adjust estimates of uninsurance; UI-fraction of the population uninsured.

Rows are:

Under-reported:	Fraction of the population under-reported (MEDS-CPS/Total Population).
Method 1:	Raw CPS Data
Method 2:	Medi-Cal from MEDS, implicitly assuming no double counting
Method 3:	Medi-Cal from MEDS, using dual coverage rate among those in the matched sample who report Medi-Cal
Method 4:	Medi-Cal from MEDS, using dual coverage rate among false negatives in the matched sample
Method 5:	Full multivariate imputation (see below)

These are the dual-coverage rates for everyone in the matched sample who reports having Medi-Cal. Figure 5.1 demonstrates that they are not the relevant population for the computation of the increase in health insurance when we correct for under-reporting. On net, our imputation moves people from the first row (does not have Medi-Cal) to the second row (has Medi-Cal). For people with other health insurance (OHI; the left column of Table 5.1), there is no net increase in health insurance/decrease in uninsurance. For people without other health insurance (the right column of Table 5.1), there is a net increase in health insurance/decrease in uninsurance.

	CPS Reports of OHI			
	No		Yes	
Medi-Cal	No	A: Uncovered	B: OHI Only	
Enrollment	Yes	C: Medi-Cal Only	D: Dual Coverage	

Figure 5.1 – Dual Coverage Rates and Adjusting Total Health Insurance Coverage (OHI: Other – non-Medi-Cal – Health Insurance)

The previous tabulations implicitly assume that those to whom we impute Medi-Cal are like those with Medi-Cal. However, this seems unlikely. We have already seen that reporting no Medi-Cal in the CPS when one actually had Medi-Cal varies with the intensity of Medi-Cal in the previous year. As the number of months of Medi-Cal enrollment drops, the probability of a false negative increases; and false negatives are the population to whom we are trying to impute Medi-Cal.

Furthermore, it seems plausible that such people are more likely to have OHI. One reason people leave Medi-Cal is that they gain private insurance. In the matched sample, we can identify such false negatives. Indeed, they have higher rates of OHI (see Table 5.1) and the difference is non-trivial. For example, for adults over the entire period, 23.9 percent of those with Medi-Cal also have OHI; for the false negatives, the figure is a quarter higher at 32.4 percent.

These tabulations are informative for three imputation methods possible with only the public-use file and aggregate tabulations from the administrative data. To understand the argument, we introduce some new notation:

$$(5.1) \quad U = T - OHI - MC + DC$$

The number of uninsured individuals (U) can be computed as the total population (T) less the count of those with other (non-Medi-Cal) health insurance (OHI), less the count of those with Medicaid insurance (MC), and adding back in those with dual coverage—other health insurance and Medicaid (DC). The results of the previous chapter imply that the CPS estimate of MC is much too small.

Given this formulation, the three imputation methods are:

Method 1: *Raw CPS Data* $U = T - OHI_{CPS} - MC_{CPS} + DC_{CPS}$:

Since Medicaid/Medi-Cal is seriously under-reported, simply using the raw data will yield an estimate of the number of uninsured that is too high.

Method 2: *Simple Administrative Data Adjustment*

$$U = T - OHI_{CPS} - MC_{CPS} + DC_{CPS} - (MC_{Admin} - MC_{CPS}):$$

Since there is significant dual coverage, estimating the number of uninsured as the total population less the CPS estimate of OHI and the administrative data estimate of Medicaid/Medi-Cal will yield an estimate of the number of uninsured that is too low.

Method 3: *Public-Use File Adjustment for Dual Coverage*

$$U = T - OHI_{CPS} - MC_{CPS} + DC_{CPS} + (1 - \delta_{CPS})(MC_{Admin} - MC_{CPS}):$$

The simple administrative data adjustment implicitly assumes no dual coverage. However, we can generate a rough estimate of dual coverage from the Public Use File, δ_{CPS} , as the fraction of those reporting Medicaid in the CPS who also report OHI.

This third Method will be appropriate if the dual-coverage rates among those reporting Medi-Cal (i.e., the union of the true positives and the false positives) equaled the dual-coverage rates among the false negatives. However, it seems plausible that it is exactly people who had private health insurance (perhaps at the end of the year) who would not report the Medi-Cal they had (perhaps at the beginning of the year). Tabulations from the matched sample are consistent with this hypothesis. (See Table 5.1). This suggests using the dual-coverage rates from the false negatives in the matched sample in the adjustment above and that doing so will yield an estimate of the number of uninsured that is larger than the basic estimate from the CPS.

This analysis is only an approximation. The full analysis has 16 cells: (matched/unmatched) x (private health insurance yes/no) x (TP, FN, FP, TN) and the total number of uninsured is:

$$U = TN[S, N] + FP[S, N] + TN[A, N] + FP[A, N]$$

where the function arguments are S/A for SSN present/SSN absent and Y/N for private health insurance/no private health insurance. Then, in both the matched and unmatched samples, there are two ways to be uninsured: (1) true negative for Medicaid and no private health insurance; or (2) false positive for Medi-Cal and no private health insurance.

Given the assumption that the MEDS data is truth for Medi-Cal and the CPS data is truth for private health insurance (we can do no better), we know the first two terms exactly and the rates from the matched sample are plausibly informative about the last two terms (i.e., rates in the unmatched sample). This analysis suggests two more estimators of the number of uninsured:

Method 4: *FN Adjustment for Dual Coverage*

$$U = T - OHI_{CPS} - MC_{CPS} - DC_{CPS} + (1 - \delta_{FN})(MC_{Admin} - MC_{CPS}):$$

Since the major concern is dual coverage among the false negatives, it seems preferable to use the rate from the matched sample's FNs, δ_{FN} . Of course, this is only possible with the matched data.

Method 5: *Full Imputation Model*: The imputation models we estimated in the previous chapter include (control for) private health insurance coverage. They thus control for dual coverage in all four cells (TP, FN, FP, TN).

Our Approach

Table 5.1 reports the results of these five adjustments for adults and children, pooled over the entire file and for the last year (the 2000 survey year, corresponding to calendar year 1999). Consider first adults over the entire panel. For this group, under-reporting is about 4.1 percent of the total. The raw CPS estimate of uninsurance is 23.5 percent. Simply adding back in the under-reporting cuts the estimated fraction of uninsured to 19.5 percent (n.b., the tables are computed with more significant digits, so the arithmetic is not exact in terms of the numbers reported in the tables).

This estimate is clearly too small. It assumes no dual coverage. Using the rate of dual coverage among those reporting Medi-Cal (Method 3) yields a slightly higher estimate of the fraction uninsured, 20.4 percent. Using either the rate of dual coverage among the false negatives (Method 4) or the full imputation model (Method 5) yields estimates of uninsurance of 20.8 percent, much lower than the simple CPS estimate (23.5 percent) and slightly higher than simply adding back in the under-reporting (19.5 percent). Thus, for adults, some adjustment for under-reporting is necessary and the reported dual coverage rate is quite close to the estimate using the matched data.

For children over the entire period, under-reporting is a much bigger problem, 9.0 percent of all children (versus 4.1 percent of all adults), and the divergence between the estimates of dual coverage is larger. Therefore, the effect of the different correction methodologies varies. The raw CPS estimate of uninsurance is 17.8 percent. Simply subtracting off the under-reporting (Method 2) cuts that estimate to 8.9 percent. The three corrections for dual coverage (Method 3, Method 4, and Method 5) successively raise the uninsurance rates. Our preferred estimate is from Method 5 (the last row). It suggests true uninsurance rates of 11.9 percent. This estimate is considerably lower than the unadjusted estimate of 17.8, but considerably higher than the no dual coverage estimate of 8.9.

The right side of the table gives the equivalent figures for the last year of our data, survey year 2000 referring to calendar year 1999. The qualitative story and the estimates of uninsurance are similar.

Revised Estimates of Uninsurance, By Demographic Group

Table 5.1 shows the effect of various corrections for dual coverage on overall estimates of uninsurance. In this section, we provide revised estimate of uninsurance using our preferred model and

the effect of the multiple imputation model. The multiple-imputation is performed for each record in the CPS. We then tabulate the multiply-imputed file by demographic characteristics—gender, single female with children, and poverty status.

Table 5.2 presents results pooled over the years 1990 to 2000. Table 5.3 presents results for 2000 only. We use the same format for both tables. The left panel refers to adults (15-65); the right panel refers to children (0-14). The rows consider subgroups related to Medi-Cal and welfare eligibility: everyone, males, females, single women with children—overall and by poverty status. For each group, we report the unadjusted uninsurance rate (the ratio of the uninsured to the population), the adjusted rate (after our imputations from the matched data), and the “Delta” — the decrease in uninsurance with imputation (the ratio of the adjusted to the unadjusted uninsurance rates minus one; note that the ratio is computed from the underlying figures with more significant digits; it thus will differ from what would be computed using the “Raw” and “Imputed” columns).

Table 5.2 has a simple story. In practice, our multiple-imputation model imputes Medi-Cal approximately in proportion to Medi-Cal coverage. This appears to be in part because dual coverage is more common among those with a lower proportion with Medi-Cal coverage. Therefore, the effect of imputation is larger the larger is the fraction of the group covered by Medi-Cal. Specifically, the effect of imputation is smallest for males, larger for females, and larger still for single females with children. Among single females with children, the effect of imputation is approximately constant among for those in poverty and much smaller for those out of poverty.

Table 5.2.
Health Insurance Coverage Rates:
Unadjusted, Adjusted, Discrepancy, Pooled Years

	Adults			Children		
	Raw	Imputed	Delta	Raw	Imputed	Delta
All	24%	21%	12%	18%	12%	33%
Male	26%	24%	7%	18%	12%	32%
Female	21%	17%	17%	18%	12%	34%
SW w/kids	27%	19%	29%	18%	12%	34%
SW w/kids <50% FPL	39%	25%	37%	30%	14%	54%
SW w/kids 50%-100% FP	30%	19%	36%	24%	13%	45%
SW w/kids 100%-150% FP	37%	25%	32%	27%	18%	33%
SW w/kids 150%-200% FP	35%	23%	36%	28%	18%	33%
SW w/kids >200% FPL	18%	15%	15%	10%	8%	14%

Note: “Raw” is the unadjusted CPS estimate; “Imputed” is the adjusted CPS estimate, based on the multiply-imputed data set; and “Delta” is the percentage (not percentage point) decrease in estimated uninsurance with imputation, as a percentage of the raw rate.

Table 5.3.
Health Insurance Coverage Rates:
Unadjusted, Adjusted, Discrepancy, 2000 Survey Year/1999 Calendar Year

	Adults			Children		
	Raw	Imputed	Delta	Raw	Imputed	Delta
All	23%	21%	12%	18%	12%	35%
Male	24%	22%	8%	17%	11%	35%
Female	22%	19%	17%	18%	12%	35%
SW w/kids	31%	22%	28%	18%	12%	35%
SW w/kids <50% FPL	39%	29%	26%	30%	14%	54%
SW w/kids 50%-100% FP	35%	24%	32%	27%	17%	38%
SW w/kids 100%-150% FP	34%	24%	29%	23%	14%	39%
SW w/kids 150%-200% FP	41%	26%	36%	28%	16%	42%
SW w/kids >200% FPL	24%	19%	24%	12%	9%	20%

Note: "Raw" is the unadjusted CPS estimate; "Imputed" is the adjusted CPS estimate, based on the multiply-imputed data set; and "Delta" is the percentage (not percentage point) decrease in estimated uninsurance with imputation, as a percentage of the raw rate.

The results for children differ slightly. As expected, there is no effect of gender (both boys and girls are covered if their family is eligible). The imputation model imputes more Medi-Cal coverage for those under poverty than for those just above poverty. In net, this makes uninsurance less strongly related to poverty status. However, it remains true that even among single women with children, uninsurance is much less common among those with incomes above twice the poverty level.

Table 5.3 presents the same results for survey year 2000/calendar year 1999. The results are similar to those for the pooled sample. The only difference is that the correction for single females out of poverty is about two-thirds larger and the correction for children in the same households is about a third larger in the later period.

Discussion

This analysis suggests that under-reporting of Medi-Cal seriously inflates our estimates of the size of the uninsured population. In addition, dual coverage is sufficiently common that ignoring it results in a significant underestimate of the size of the uninsured population even after we correct for under-reporting.

Thus, some correction for dual coverage is needed. On a priori grounds, this dual-coverage estimate from the false negatives in the matched sample seems to be a preferable estimate of the unmatched false negatives than the simple CPS public-use file estimate. Evidence from the matched sample suggests that dual coverage is slightly more common among false negatives and also more common than the simple CPS public-use file estimate. Using this plausibly better estimate yields a slightly higher estimate of the uninsured. Using the full imputation model yields an estimate that is slightly higher. The size of the corrections will vary with the magnitude of the under-reporting and the amount of dual coverage.

Our best estimates of dual coverage suggest that the raw CPS figures over-estimate uninsurance by about three percentage points for adults and eight percentage points for children. These are sizable over-estimates. They imply that the problems of lack of health insurance are non-trivial, but that they are considerably smaller than what would be implied by the simple CPS tabulations.

6. New Estimates of Medi-Cal and Welfare Enrollment Rates

If the only question of interest was: “How many people are enrolled in welfare/Medicaid?”, we could answer that question directly from the administrative data. However, both researchers’ and policy makers’ interest typically goes well beyond the number of people enrolled to concerns about take-up rates. The question of interest is what share of the target population is actually enrolled in the program of interest. This is a rate that cannot be measured with administrative data. While the numerator (i.e., the number of people enrolled) is available in the administrative data, the denominator (i.e., the number of people in the target population) is not, because the administrative data only includes information on those actually enrolled in the program.

Therefore, to estimate take-up rates, analysts generally turn to survey data for both pieces of information (i.e., the number of people in the target population and the number of people enrolled). Unfortunately, we have seen that actual enrollment is seriously under-reported in survey data, so take-up rates based on these data will also be under-reported. Furthermore, we have seen that under-reporting is not random. Non-reporting is more common among those who are covered for less of the year and closer to the border of eligibility.

In this chapter, we use the adjusted California CPS data based on the analyses and methods described in the previous chapters to generate new estimates of what we call “enrollment rates”. We note that the estimates here are not pure take-up rates. Conventional pure take-up rate estimates attempt to impute eligibility for Medi-Cal based on all of the survey information and Medi-Cal program rules. Here, we simply compute enrollment given conventional demographic variables and income relative to the poverty line. We do not attempt a full eligibility simulation.

Pooled Results

We begin by pooling across all of the years in our analysis, 1990-2000. Table 6.1 presents our basic results for Medi-Cal, in the format used for all the tables that follow. The left panel refers to adults (15-65); the right panel refers to children (0-14). The rows consider subgroups related to Medi-Cal and welfare eligibility: everyone, males, females, single women with children—overall and by poverty status. For each group, we report the unadjusted CPS enrollment rate (the ratio of enrollees to the population), the adjusted rate (after our imputations from the matched data), and the “Delta”—the increase with imputation (the ratio of the adjusted to the unadjusted enrollment rates minus one; note that the ratio is computed from the underlying figures with more significant digits; it thus will differ from what would be computed using the “Raw” and “Imputed” columns).

Table 6.1.
Enrollment Rates: Unadjusted, Adjusted, Discrepancy
Medi-Cal, Pooled Years

	Adults			Children		
	Raw	Imputed	Delta	Raw	Imputed	Delta
All	10%	14%	42%	27%	36%	34%
Male	7%	10%	38%	26%	35%	33%
Female	12%	18%	45%	27%	36%	34%
SW w/kids	28%	39%	41%	27%	36%	34%
SW w/kids <50% FPL	48%	65%	34%	59%	77%	31%
SW w/kids 50%-100% FP	60%	74%	23%	64%	78%	22%
SW w/kids 100%-150% FP	39%	55%	40%	40%	53%	34%
SW w/kids 150%-200% FP	22%	40%	86%	21%	36%	72%
SW w/kids >200% FPL	8%	14%	77%	7%	10%	55%

Note: "Raw" is the unadjusted CPS estimate; "Imputed" is the adjusted CPS estimate, based on the multiply-imputed data set; and "Delta" is the percentage (not percentage point) increase in estimated enrollment with imputation.

Overall Medi-Cal enrollment increases by about 40 percent when we adjust for under-reporting using our imputation model. The increases are slightly larger for adults (42 percent) and slightly smaller for children (34 percent). Consistent with an explanation of under-reporting due to stigma, the increases are smallest for those in poverty, largest for those between one and two times poverty, and large for those at more than twice poverty.

We note the anomalous result that, even after adjustment, adults (but to a much lesser degree children) at less than half the poverty line have lower enrollment rates than those between half the poverty line and the poverty line. Welfare is considered income in the computation of the poverty line. In California, welfare takes a family to more than half the poverty line. Therefore, the families at less than half the poverty line are unlikely to have contact with the welfare system, including Medi-Cal.

Table 6.2 presents comparable estimates for welfare. Consistent with earlier results, compared to Medi-Cal, the levels of enrollment are lower and the adjustments have a larger effect. The average adjustment more than doubles enrollment rates. The adjustments are similar across children and adults. The adjustments are smaller for those near poverty, and larger for those out of poverty.

Table 6.2.
Enrollment Rates: Unadjusted, Adjusted, Discrepancy
Welfare, Pooled Years

	Adults			Children		
	Raw	Imputed	Delta	Raw	Imputed	Delta
All	3%	6%	108%	11%	23%	105%
Male	1%	4%	349%	11%	23%	108%
Female	5%	9%	69%	12%	23%	101%
SW w/kids	15%	24%	58%	12%	23%	101%
SW w/kids <50% FPL	25%	42%	66%	26%	59%	131%
SW w/kids 50%-100% FP	41%	55%	34%	35%	58%	67%
SW w/kids 100%-150% FP	20%	33%	61%	15%	30%	105%
SW w/kids 150%-200% FP	9%	19%	115%	6%	17%	203%
SW w/kids >200% FPL	3%	6%	135%	1%	4%	206%

Table 6.3 and Table 6.4 present the equivalent results for the last year of our data, the 2000 survey referring to the 1999 calendar year. Consistent with our basic analysis of under-reporting in Chapter 2, the results are similar for Medi-Cal. However, for welfare, there has been a sharp increase in non-reporting. Correspondingly, the deltas are much higher in 2000 than in the pooled sample, both overall for adults (169 percent versus 109 percent) and for children (263 percent versus 132 percent).

Table 6.3.
Enrollment Rates: Unadjusted, Adjusted, Discrepancy
Medi-Cal, 2000 Survey/1999 Calendar Year

	Adults			Children		
	Raw	Imputed	Delta	Raw	Imputed	Delta
All	9%	13%	42%	26%	35%	38%
Male	7%	9%	37%	27%	36%	36%
Female	11%	17%	45%	25%	35%	40%
SW w/kids	25%	37%	46%	25%	35%	40%
SW w/kids <50% FPL	49%	59%	21%	58%	75%	30%
SW w/kids 50%-100% FP	48%	63%	31%	57%	70%	22%
SW w/kids 100%-150% FP	41%	55%	34%	40%	55%	39%
SW w/kids 150%-200% FP	24%	45%	84%	26%	47%	77%
SW w/kids >200% FPL	8%	16%	104%	8%	12%	66%

Table 6.4.
Enrollment Rates: Unadjusted, Adjusted, Discrepancy
Welfare, 2000 Survey/1999 Calendar Year

	Adults			Children		
	Raw	Imputed	Delta	Raw	Imputed	Delta
All	2%	4%	169%	5%	18%	233%
Male	0%	3%	524%	6%	18%	215%
Female	3%	6%	117%	5%	18%	253%
SW w/kids	8%	18%	123%	5%	18%	253%
SW w/kids <50% FPL	19%	35%	81%	17%	58%	242%
SW w/kids 50%-100% FP	17%	37%	114%	11%	41%	285%
SW w/kids 100%-150% FP	17%	31%	82%	11%	27%	147%
SW w/kids 150%-200% FP	5%	16%	193%	6%	21%	230%
SW w/kids >200% FPL	1%	6%	482%	0%	3%	3549%

Discussion

This chapter has considered the effect of correcting for under-reporting of program enrollment in sub-populations using our imputation model for program enrollment rates. The earlier chapters found substantial under-reporting. That under-reporting yields large under-estimates of enrollment rates. Furthermore, the variation in enrollment rates is not simple. Overall, the corrections appear to be smallest for those in deepest poverty and larger for those with only borderline eligibility. This pattern is consistent with greater stigma for the borderline eligible. It is also consistent with the borderline-eligible only being enrolled for part of the year; we have seen that those enrolled for part of the year are less likely to report enrollment in the CPS.

7. Conclusion

This report describes analyses of matched CPS-MEDS data for California. The CPS data are known to have substantial under-reporting of Medi-Cal enrollment and even larger under-reporting of welfare enrollment. The matched data—along with some auxiliary assumptions—generate adjusted estimates of who is covered by Medi-Cal in the CPS and of true health insurance coverage rates. In brief, we find that adjusting substantially cuts the estimates of the uninsured population and substantially increases estimates of enrollment rates. Given that these results confirm that the CPS data significantly undercounts enrollment, do these results suggest any solution? Happily the answer appears to be “yes.” Non-reporting is differential, but the differentials are second-order compared to the non-reporting itself. Simple ratio adjustments with a simple correction for dual coverage (e.g., from those in the CPS who report Medi-Cal coverage) are likely to eliminate most of the bias.

Such simple ratio adjustments can be computed from unmatched tabulations from the CPS and the MEDS. These under-count rates do vary over time. Thus, current official tabulations are needed. CDHS and CMS already publish some such tabulations. To correct the CPS, ideally one would use tabulations slightly different from those currently published. The necessary tabulations would consider any receipt in the past year, with consistent breaks by program and age. CDHS could easily generate the requisite tabulations.

When considering these results, it is important to note that the matched data are only for California. It seems likely that California is very different from the rest of the nation. First, relative to the most of the nation, it has high income eligibility thresholds for welfare and Medicaid and high enrollment rates in both programs. However, the state has a large immigrant population, much of it undocumented. Partially as a result, the state has high rates of uninsurance. It seems plausible that the high rates of uninsurance imply that false negatives to the Medi-Cal questions are more likely to be truly uninsured, while in other states false negatives on the Medi-Cal question might be more likely to have health insurance. This question could be explored in further research by considering how dual coverage rates among those who report having Medicaid varies across states.

Appendix A. Detailed Notes on File Construction and Matching

This appendix describes in detail the procedures we used to create our analysis files. We begin by describing the raw data files we received, how we merged them together, and the results of our efforts to eliminate false matches. This appendix concludes with a discussion of the congruence of race/ethnicity coding across the two data sources.

A.1. The Raw Data

This project's analyses are made possible by the availability of scrambled SSNs (referred to by Census as PIKs) on both the CPS data and on the MEDS data. Specifically, we received three files:

- *MEDS File*: We received an extract from the Medi-Cal Eligibility Data System (MEDS). The MEDS is the official roster of those on Medi-Cal in California. Conceptually the file contained one record for each person who appeared in the MEDS data in any month between January 1987 and December 2002.²³ Given its purpose, anyone who was covered by Medi-Cal during this period (including welfare/AFDC/TANF/CalWORKs) should have a record in the file. The record contains basic demographics (gender, date of birth, race/ethnicity, language) and for each month January 1989 to December 2000²⁴ the aid code (i.e., type of Medi-Cal coverage), eligibility status, case id, county-of-residence, and zip code. Finally, the file contained a Census generated PIK (Protected Identification Key). This file contained 22,848,715 persons.
- *CPS File*: We created the CPS file directly from the March Public Use Files. This file contained 1,588,115 observations from 1990 to 2000. The CPS rotation group structure implies that most people will appear in two successive CPS files. This observation count thus counts such people twice. We have made no correction for that correlation.
- *Cross-walk File*: For every person in the March CPS for interview years 1990 to 2000, the cross-walk file contained the PIK (Protected Identification Key), the CPS Household Number, the CPS Person Number. The file contained 1,347,282 observations (including interviews inside and outside of California).

²³ In fact, to minimize file size, the MEDS data arrived in three files. To understand the three files it is useful to note that the underlying MEDS file was created from archived version of the MEDS created every six months (the "December cut" and the "June cut"), each with 15 months of history. The first file had the time invariant data (gender, date of birth, race/ethnicity, language), taken from the most recent file. The second file had the invariant information from an underlying MEDS "cut file" (county-of-residence, case id, zip code), one record for each person, for each "cut file." The third file had the information that varied every month. The MEDS "cut files" are overlapping. We used the information from the most recent file.

²⁴ Note that we have monthly data for a narrower window (January 1999 to December 2000) than the period over which one would have needed to have been covered by Medi-Cal to be in the file at all (January 1997 to December 2002). The net result is that people who were covered by Medi-Cal in the broader window, but not in the narrower window, will have MEDS records that never show Medi-Cal coverage (all "zeros"). See the next section for a discussion of the implications of this distinction.

Unfortunately, the Cross-walk Files we received were not consistent across years. For every individual 15 and older, the CPS interviewer requests a Social Security Number (SSN). In practice, not everyone supplies an SSN; and not all of those SSNs that are reported are correct (some due to simple errors of memory or transcription; some due to deliberate obfuscation). As far as we have been able to ascertain, over the CPS years we are analyzing (1990-2000), there have been no significant changes in the procedures for the collection of SSNs, nor for the verification of SSNs.

We note that the CPS files we use are of two types. Some of the files are “unvalidated.” This apparently means that the SSNs from which the PIKs were constructed are as recorded by the Census interviewer.

This is in contrast to “validated” files. “Validated” files differ in two ways. First, those SSNs that were provided were checked against Social Security Administration SSN records (apparently checking for a match on name, gender, and birth date). When the CPS provided information that did not match the information in the SSA records for that SSN, the SSN was dropped (even though an SSN had been provided).

Second, those individuals who did not or could not provide an SSN were asked for permission to use SSA files to impute an SSN. When permission was granted, the name, gender, and birth date provided were matched against SSA records. When a match was found, that SSN was appended to the record (even though no SSN had been provided at the interview).

As we note in the body of the report. In validated files, a higher percentage of records have SSNs and a lower fraction of the matches are “bad.”

We also note that our analysis would be easier and of higher quality if we had a consistent time series (ideally, all validated; but alternatively all unvalidated). Presumably, unvalidated versions of the files for each year once existed. We were, however, informed that unvalidated versions of the validated files were not available; and that, Census could not provide validated versions of the currently unvalidated data within the time frame of the project (i.e., in six to nine months).

A.2. Matching the CPS and the MEDS

In principle, the cross-walk file should allow us to link MEDS records to CPS individuals. The link, however, will never be complete.

For example, in order to complete the link, we must have an SSN for each CPS individual. However, as noted above, not everyone has an SSN in the files we received. Individuals under 15 years of age were not asked to provide SSNs.²⁵ Many of those asked to provide SSNs were unable to do so or refused. Only some of those asked gave permission for SSA to impute an SSN, and only for some of those people who gave permission, did SSA actually provide an SSN.

A.3. Verifying Matches

The previous discussion concerns mechanical matches of records in the two files. As we discuss below in detail, there must be false non-matches in the matching process: People who were actually

²⁵ Given the CPS rolling panel structure, it might be possible to recover SSNs for half of the 14 year olds, by using conventional CPS matching methods (to match the SSN provided at their second March interview (when they were 15) to their first March interview (when they were 14)). We have not done so.

enrolled in Medi-Cal, but for whom we cannot match MEDS and CPS records. Given the non-trivial rates of not providing an SSN, it seems likely that most of these false non-matches are people who refused to provide an SSN. In addition, some people are likely to have provided an incorrect SSN that did not match to any MEDS record (and which was not caught by the SSA validation process, perhaps because this was a year in which no validation was done). There does not appear to be anything that we can do directly about this. However, we discuss below our indirect adjustments for this issue.

There is also the possibility of false matches: People who provided an SSN that matched to the MEDS, but for whom the two records do not represent the same person. In unvalidated years, this might represent people who gave the wrong SSN (a memory error, a transcription error, or deliberate obfuscation). In validated years, such errors should almost always have been caught by comparisons of name, gender, and age between the CPS responses and the SSA files.

To assure that both the CPS and MEDS records with the same SSN truly referred to the same individual, we verified the correspondence of gender and age between the MEDS data and the CPS data (where we use the term “verified” to refer to our cross-check and the term “validated” to refer to SSA’s cross-check). We considered matches to be verified if gender matched and age differed by no more than one year (i.e., we did not require an *exact* match on age). Matches not meeting both of these criteria were deemed false or bad matches and dropped at this verification stage.

In such cases where the SSNs match, some (perhaps most) of the failure to “verify” is probably caused by incorrect recording of SSN, gender, or age in the MEDS or in the CPS. Note that in validated years, the validation process should have caused the SSN to be dropped from CPS records for which the SSN, gender, and age information did not match the SSA administrative data. The net result should have been higher quality SSN, gender, and age information in the CPS records with validated SSNs.

Table A.1 provides some additional information on the results of the verification. The table stratifies by validated and unvalidated years. For each set of years, we report the distribution of correspondence by gender and age. Note that the match rates improve with validation. For the validated data, we accept 95 percent of the matches; while for the unvalidated data, we accept 92 percent of the matches. The matches deemed rejected are entered in italics.

Table A.1.
Age Differential (CPS Age - MEDS Age)
(Percent within Validation Status)

Gender Match	Validated		Unvalidated	
	Yes	No	Yes	NO
--	1%	0%	2%	2%
-1	2%		4%	
0	91%	1%	84%	1%
+1	2%		4%	
++	2%	1%	3%	0%

Note: Rows are age difference (CPS Age - MEDS Age). CPS age is as per response at survey; MEDS age computed based on birth date and CPS interview date.

“--” MEDS age is greater than CPS age by two or more years; “-1” MEDS age is greater than CPS age by exactly one year; “0” MEDS age equals CPS age; “+1” CPS age is older than MEDS age by exactly one year; “++” CPS age is older than MEDS age by more than two or more years for gender match and discrepancy for no gender match.

Appendix B. Regression Specification for Response Errors

Our approach to extrapolating from the matched sample to the unmatched sample requires imputing the probability of (what we call in Chapter 3) “imputational false negative/positive rates” to each individual. In addition, the discussion in Chapter 3 emphasizes that the corresponding “behavioral false negative/positive rates” are of substantive interest; they describe how the probability of response errors varies with other recorded survey responses. We approximate those rates using logistic regression models. This appendix describes our strategy and the details of our specification.

B.1. The Basic Approach

We model eight types of reporting errors, defined by all combinations of three binary criteria. We model two outcomes: Any Medi-Cal and welfare. For each outcome, we model four response errors: False negatives and false positives, according to the “imputational” and “behavioral” discussed in Chapter 3. The imputational model gives the probability of an incorrect response, given the recorded survey response. The behavioral model gives the probability of an incorrect given the true receipt status. Recall that we treat the administrative data as “truth”, so “incorrect” is defined relative to the administrative data.

Corresponding to the binary nature of the outcomes (a correct report or an incorrect response) and our need for predictions in the unit interval, we estimate the models by logistic regression. Estimation proceeds ignoring the weights. We then adjust the standard errors for the weights using a Huber/linearization approximation. Estimation is done in SAS using a PROC IML routine provided by Dan McCaffrey of the RAND Statistics Group.

The final logistic regression specifications whose results were reported below were generated by a simple regression pruning procedure based on three sets of variables.

- *Pure Time Effects:* We have eleven years of data, so a fully specified model would include a constant and ten year effects. To allow pruning of the specification, we specify the year effects as a linear time trend, an offset for the unvalidated years, a linear time trend specific to the unvalidated years, and year dummies for 1992, 1993, 1994, 1995, 1997, 1998, and 2000. The excluded years are 1990, 1991, 1996, and 1999.
- *Time Invariant Covariate Effects:* We define a vector of covariates related to plausible eligibility for Medi-Cal/Welfare. See below for the exact definitions.
- *Linear Covariate \times Time Interactions:* We allow the effect of each of the covariates to vary linearly (with survey year) over the eleven years of data in our analysis file.

Given this list of potential covariates, we proceed in four steps:

1. We begin with all of the pure time effects and the time invariance covariate effects included.
2. We then delete all variables (from the time effects and from the time invariant covariate effects) that do not have a p-value less than 0.10. Note that we only do this once. We do not iterate. We also do not require that the p-values in the final model be less than 0.10.

3. We then include linear time interactions with all of the time invariant covariate effects that remained in the model after the deletion in the second step.
4. We then again delete all interactions that are not significantly different from zero at $p=0.10$. This is our final model. We report these results below.

The candidate time invariant covariates are as follows:

- *Age Dummies*: For 15-25, 36-65, and 46-65. The excluded category is 26-35. The dummies are deliberately overlapping to allow the variable deletion to proceed. Younger people have lower earnings potential and are thus more likely to be eligible for the program.
- *Gender Effects*: A dummy variable for male (female is the excluded category) and an identical set of age dummies interacted with the male dummy. Since children usually reside with the mother, it is difficult for males to qualify for welfare and therefore for Medi-Cal.
- *Race/Ethnicity*: Dummy variables for black and Hispanic. The categories are not tied; i.e., a black-Hispanic is logically possible, though not common. Blacks and Hispanics have lower earnings potential and are thus more likely to be eligible for the program.
- *Education*: Dummy variables for high school drop out and at least some college (including college graduates). Less educated individuals have lower earnings potential and are thus more likely to be eligible for the programs.
- *Poverty*: Dummy variables for income less than half the poverty line, less than the poverty line, less than one and a half times the poverty line, and less than twice the poverty line. Again, note that the categories are deliberately overlapping. Program eligibility is explicitly conditioned on income.
- *Family Structure*: Dummy variables for any children in the household and a single female-headed household with children. The presence of children is required to qualify for AFDC-UP/two parent welfare.
- *Other Health Insurance*: A dummy variable for the presence of health insurance other than Medi-Cal (primarily employer sponsored health insurance, but also other private health insurance, and other public health insurance—Medicare and CHAMPUS/VA). Someone with other health insurance (during part of the year) might be less likely to report Medi-Cal (or be rich enough to be unlikely to have received welfare). Note that the health insurance questions are asked separately from the welfare and Medicaid/Medi-Cal questions that are asked in the context of program participation.
- *Program Participation*: When the result is not logically implied by the response (e.g., everyone with welfare has Medi-Cal), we include dummy variables for welfare or Medi-Cal-only (i.e., Medi-Cal, but not welfare). These variables appear in none of the welfare models. The welfare variable appears in both Medi-Cal false positive models; i.e., someone who also reports welfare in the survey might be less likely to falsely report having Medi-Cal.

The actual list of regressors at the first step is smaller than is implied by this list. The Census Bureau's confidentiality rules prevent disclosure of regression results that include dummy variables with insufficient number of observations in either value of the binary outcome variable. The exact minimum cell counts are themselves confidential. In order to assure that we could release the results, variables not meeting the standard were not included in the initial step.

B.2 Detailed Logistic Regression Results

Table B.1 and Table B.2 report the final logistic regression models for Medi-Cal and welfare respectively. Each table reports four regressions: Behavior/Imputational x False Negative/False Positive. For each model, the table reports the logistic regression coefficient and its standard error. The standard errors are computed using sample weights and robust standard errors. Empty cells indicate that the variable was not included in the final model. No row is included for variables that never entered any of the four final models.

Table B.1.
Detailed Logistic Regression Results for Medi-Cal

	Behavioral				Imputational			
	False Negative		False Positive		False Negative		False Positive	
	Beta	S.E.	Beta	S.E.	Beta	S.E.	Beta	S.E.
INTERCEPT	-0.626	0.129	-4.384	0.156	-3.111	0.129	-0.973	0.119
YEAR_DEV			-0.073	0.017	0.088	0.013		
UNVALID	-0.378	0.079	0.515	0.182				
UNVALID_YR			0.071	0.029				
Y1994			0.169	0.112				
Y1997			0.267	0.124				
Y1998			0.153	0.141				
A1525					0.214	0.103		
A3665			-0.183	0.085	-0.340	0.095		
A4665	-0.445	0.105	0.369	0.103	-0.382	0.137		0.161
M1525					0.872	0.168	-0.648	0.176
M3665					0.801	0.160	-0.626	0.174
M4665							0.229	0.137
MALE	-0.258	0.084	-0.202	0.077	-1.278	0.129	0.530	0.078
HISPANIC	0.707	0.110					0.819	0.138
BLACK	0.367	0.133	0.927	0.122	1.012	0.109	0.431	0.145
HSDO			0.183	0.083				
SOCO	0.504	0.129	-0.168	0.146	-0.532	0.071	0.545	0.108
POV_LT_05	0.366	0.176					0.718	0.118
POV_LT_10	-0.374	0.094	0.579	0.089	0.220	0.083		
POV_LT_15	-0.708	0.112	0.439	0.117	0.169	0.095	-0.429	0.150
POV_LT_20	-0.623	0.147	0.900	0.118	1.169	0.093	-0.937	
KIDSINHH	0.776	0.147	0.531	0.085	0.876	0.079		0.148
SFWKIDS	-0.805	0.131	0.625	0.094	0.352	0.091	-0.529	0.086
OHI	0.819	0.078	-1.167	0.083	-1.094	0.102	0.379	0.088
C_W							-0.414	0.027
A4665_YR	0.027	0.021			-0.045	0.025		
SOCO_YR	0.036	0.026					0.048	0.016
POV_LT_05_YR	-0.095	0.035	0.078	0.026				
POV_LT_20_YR	-0.073	0.024					-0.096	0.024
KIDSINHH_YR	0.065	0.025						
SFWKIDS_YR	-0.025	0.024					-0.036	
OHI_YR					-0.027	0.018		
N		5,961		50,339		50,751		5,549

Table B.2.
Detailed Logistic Regression Results for Welfare

	Behavioral				Imputational			
	False Negative		False Positive		False Negative		False Positive	
	Beta	S.E.	Beta	S.E.	Beta	S.E.	Beta	S.E.
INTERCEPT	0.866	0.294	-7.658	0.318	-5.253	0.164	-1.524	0.238
YEAR_DEV	0.089	0.029	-0.061	0.024				
UNVALID	-0.696	0.220	0.951	0.275	0.018	0.012	-0.467	0.157
UNVALID_YR	-0.104	0.038	0.102	0.046				
A1525	0.530	0.174	-0.953	0.149				
A3665	0.206	0.122						
A4665	-0.230	0.276			-1.127	0.146	0.650	0.352
M1525					0.533	0.106		
M4665					0.782	0.207		
MALE	0.906	0.137			-0.485	0.098		
HISPANIC	0.347	0.091			-0.344	0.096	0.933	0.132
BLACK					0.288	0.153		
HSDO	0.371	0.095	0.364	0.122	0.312	0.084	0.297	0.134
SOCO					-0.527	0.102		
POV_LT_05	0.624	0.214					0.835	0.168
POV_LT_10	-0.046	0.174	0.755	0.147	0.438	0.087		
POV_LT_15			0.503	0.215	0.200	0.115		
POV_LT_20	-0.642	0.166	1.076	0.251	0.973	0.121	-0.687	0.221
KIDSINHH	-0.520	0.232	1.757	0.265	1.709	0.121	0.668	0.178
SFWKIDS	-0.726	0.119	1.922	0.125	0.606	0.091	-0.135	0.063
OHI	1.076	0.131	-1.426	0.164	-0.845	0.088		
C_O					2.017	0.100		
C_O_YR					-0.040	0.018		
HISPANIC_YR					0.053	0.016		
A1525_YR	-0.088	0.032						
A4665_YR	-0.139	0.052						
POV_LT_05_YR	-0.098	0.045						
POV_LT_10_YR	0.084	0.032						
N		3,159		54,141		54,414		1,886

Bibliography

- Alpha Center. "The Impact of Federal Welfare Reform on Medicaid: Medicaid Enrollment Declines as De-linking Yields Eligibility Confusion." State Initiative in Health Care Reform. R30, pp. 4-10, Washington, DC, 1999.
- Alpha Center. A Survey of Surveys: What Does it Take to Obtain Accurate Estimates of the Uninsured? *State Coverage Initiatives*, No. 1, March 2000.
- Bavier, R. (1999). *An Early Look at the Effects of Welfare Reform*. Unpublished manuscript, Office of Management and Budget, Washington, DC.
- Bavier R. (2003). *Non-Economic Factors in Early Welfare Caseload Declines*. Unpublished manuscript, Office of Management and Budget, Washington, DC.
- Beauregard, Karen M., Susan K. Drilea, and Jessica P. Vistnes. "The Uninsured in America-1996." *MEPS Highlights*, No. 1, May 1997.
- Bennefield, R. L. "Health Insurance Coverage: 1995." *Current Population Reports*, pp. 60-199. Washington, DC: U.S. Bureau of the Census, 1996.
- Bennefield, Robert L. "A Comparative Analysis of Health Insurance Coverage Estimates: Data from CPS and SIPP." Proceedings from the 1996 Joint Statistical Meetings, American Statistical Association, Chicago, August 1996c.
- Bennefield, Robert L. "Dynamics of Economic Well-Being, Health Insurance, 1993 to 1995." *Current Population Reports*, Washington, DC: U.S. Department of Commerce, Economics and Statistics Administration, July 1998, pp. 70-64.
- Bennefield, Robert L. "Who Loses Coverage and for How Long?" *Current Population Reports*, Washington, DC: Census Bureau, May 1996b, pp. 70-54.
- Bennefield, Robert, L. "Health Insurance Coverage 1995." *Current Population Reports*, Washington, DC: Census Bureau, pp. 60-195, 1996a.
- Bilheimer, Linda T. "CBO Testimony on Proposals to Expand Health Coverage for Children." Testimony before the Subcommittee on Health, U.S. House of Representatives, Committee on Ways and Means, Washington, DC, April 8, 1997.
- Blumberg, S. J. and M. L. Cynamon. "Misreporting Medicaid Enrollment: Results of Three Studies Linking Telephone Surveys to State Administrative Records." Paper delivered at the 7th Conference on Health, Survey Research Methods, Williamsburg, VA, 1999.
- Bradburn, N. M., L. J. Rips, and S. K. Shevell. "Answering Autobiographical Questions: The Impact of Memory and Inference on Surveys." *Science*, Vol. 236, No. 4798, pp. 157-161, 1987.
- Call, K. T., A. S. Somrners, R. Feldman, T. Rockwood, Y. Jonk, and B. Dowd. Minnesota Health Access Survey, 1999 Final Report. University of Minnesota, School of Public Health, Division of Health Services Research and Policy, Minneapolis, MN, 1999.

- Call, Kathleen Thiede, Gestur Davidson, Anna Stauber Sommers, Roger Feldman, Paul Farseth, and Todd Rockwood. "Uncovering the Missing Medicaid Cases and Assessing their Bias for Estimates of the Uninsured." *Inquiry*, Vol. 38, No. 4, Winter 2001/2002, pp. 396-408.
- Card, David, Andrew K. G. Hildreth, and Lara D. Shore-Sheppard. "The Measurement of Medicaid Coverage in the SIPP: Evidence from California, 1990-1996." Cambridge, MA: National Bureau of Economic Research (NBER), Working Paper No. 8514, October 2001. [<http://papers.nber.org/papers/w8514.pdf>]
- Congressional Budget Office. *How Many People Lack Health Insurance and For How Long*, 2003.
- Cutler, David and Jonathan Gruber. "Does Public Insurance Crowd out Private Insurance?" *Quarterly Journal of Economics*, CXI, 1996, pp. 391-430.
- Cutler, David and Jonathan Gruber. "Medicaid and Private Insurance: Evidence and Implications." *Health Affairs*, Vol. 17, No. 1, 1997, pp. 194-200.
- Dubay L. and G. Kenney. "Lessons from the Medicaid Expansions for Children and Pregnant Women: Implications for Current Policy." Statement for Hearing on Children's Access to Health Coverage, Subcommittee on Health, U.S. House Committee on Ways and Means, April 8, 1997.
- Dubay, L., and G. Kenney. "Effects of Medicaid Expansions on Insurance Coverage of Children." *Future of Children*, Vol. 6, No. 1, pp. 152-161, 1996.
- Fay, R. E. "An Analysis of Within-Household Undercoverage in the Current Population Survey." Paper delivered at the U.S. Bureau of the Census Annual Research Conference, Washington, DC, 1989.
- Fronstin, P. "Expanding Health Insurance for Children: Examining the Alternatives." Washington, DC: Employee Benefit Research Institute, 1997a.
- Fronstin, P. "Trends in Health Insurance Coverage." Washington, DC: Employee Benefit Research Institute, 1997b.
- Fronstin, Paul and Rachel Christensen. "The Relationship Between Income and the Uninsured." *EBRI Notes*, No. 3, Employee Benefit Research Institute, March 2000, pp. 1-4.
- Fronstin, Paul. "Sources of Health Insurance and Characteristics of the Uninsured: Analysis of the March 1996 Current Population Survey." *EBRI Issue Brief*, No. 179, Washington, DC: EBRI, November 1996.
- Fronstin, Paul. "Trends in Health Insurance Coverage." *EBRI Issue Brief*, No. 185, Washington, DC: EBRI, May 1997.
- Giannarelli, L. *An Analyst's Guide to TRIM2: The Transfer Income Model, Version 2*. Washington, DC: The Urban Institute, 1992.
- Groves, R. M. *Survey Errors and Survey Costs*. New York, NY: Wiley and Sons, 1989.
- Gruber, Jonathan. "Medicaid." Cambridge, MA: National Bureau of Economic Research (NBER), Working Paper 7829, August 2000. [<http://papers.nber.org/papers/w7829.pdf>]

- Hainer, P., C. Hines, E. Martin, and G. Shapiro. "Research on Improving Coverage in Household Surveys." Paper delivered at the U.S. Bureau of the Census Annual Research Conference, Washington, DC, 1988.
- Hall, J., G. Kenney, G. Shapiro, and I. Flores-Cervantes. "Bias from Excluding Households without Telephones in Random Digit Dialing Surveys: Results of Two Surveys." *Proceedings of the Survey Research Methods Section of the American Statistical Association*, 1999, pp. 382-387.
- Hogan, H. "The 1990 Post-Enumeration Survey: Operations and Results." *Journal of the American Statistical Association*, Vol. 88, No. 423, pp. 1047-1060, 1993.
- Holahan, J., C. Winterbottom, and S. Rajan. "A Shifting Picture of Health Insurance Coverage." *Health Affairs*, Vol. 14, No. 4, pp. 253-264, 1995.
- Kronick, R. "Health Insurance, 1979-1989: The Frayed Connection between Employment and Insurance." *Inquiry*, Vol. 28, 1991, pp. 318-332.
- Krosnick, J. A., S. Narayan, and W. R. Smith. "Satisficing in Surveys: Initial Evidence." *New Directions for Program Evaluation*, Vol. 70, pp. 29-44, 1996.
- Levit, Katharine R., Gary L. Olin, and Suzanne W. Letsch. "Americans' Health Insurance Coverage, 1980-91." *Health Care Financing Review*, Vol. 14, No. 1, Fall 1992, pp. 31-57.
- Lewis, K., M. Ellwood, and J. L. Czajka. *Counting the Uninsured: A Review of the Literature*. Washington, DC: The Urban Institute, 1998.
- Lewis, Kimball, Marilyn Ellwood, and John L. Czajka. "Counting the Uninsured: A Review of the Literature." Washington, DC: Urban Institute, Assessing the New Federalism, Occasional Paper No. 8, July 1998.
- Long, Stephen H. and M. Susan Marquis. "Some Pitfalls in Making Cost Estimates of State Health Insurance Coverage Expansions." *Inquiry*, Vol. 33, Spring 1996, pp. 85-91.
- Marquis, K. and J. Moore. "Measurement Errors in the Survey of Income and Program Participation (SIPP) Program Reports." Paper read at 1990 Annual Research Conference, August 1990.
- Martini, A. "Seam Effect, Recall Bias, and the Estimation of Labor Force Transition Rates from SIPP." Paper read at American Statistical Association, Section on Survey Research Methods, August 1989.
- Moyer, M. Eugene. "A Revised Look at the Number of Uninsured Americans." *Health Affairs*, Summer 1989, pp. 102-110.
- Nadeau, R. and R. G. Niemi. "Educated Guesses: The Process of Answering Factual Knowledge Questions in Surveys." *The Public Opinion Quarterly*, Vol. 59, No. 3, pp. 323-346, 1995.
- Nelson, Charles T. and Robert J. Mills. "The March CPS Health Insurance Verification Question and Its Effect on Estimates of the Uninsured." U.S. Bureau of the Census, Housing and Household Economic Statistics Division, August 2001. [<http://www.census.gov/hhes/hlthins/verif.html>]
- Perry, M., S. Kannel, R. B. Valdez, and C. Chang. "Medicaid and Children: Overcoming Barriers to Enrollment, Findings from a National Survey." Washington, DC: Kaiser Commission on Medicaid and the Uninsured, 2000.

- Presser, S. "Is Inaccuracy on Factual Survey Items Item-Specific or Respondent-Specific?" *The Public Opinion Quarterly*, Vol. 48, No. 1B, pp. 344-355, 1984.
- Rajan, Shruti, Stephen Zuckerman and Niall Brennan. "Confirming Insurance Coverage in a Telephone Survey: Evidence from the National Survey of America's Families." *Inquiry*, Vol. 37, No. 3, Fall 2000, pp. 317-327.
- Schuman, H. and S. Presser. *Questions and Answers in Attitude Surveys: Experiments on Question Form, Wording, and Context*. New York, NY: Academic Press, 1981.
- Selden, T. M., J. S. Banthin, and J. W. Cohen. "Medicaid's Problem Children: Eligible but not Enrolled." *Health Affairs*, Vol. 17, No. 3, pp. 192-200, 1998.
- Shapiro, G., G. Diffendal, and D. Cantor. "Survey Undercoverage: Major Causes and New Estimates for Magnitude." Paper delivered at the U.S. Bureau of the Census Annual Research Conference, Washington, DC, 1993.
- Short, Pamela Farley. "Counting and Characterizing the Uninsured." ERIU Working Paper 2, 2001. [<http://www.umich.edu/~eiru/pdf/wp2.pdf>]
- StatCorp. *Stata Statistical Software: Release 5.0*. College Station, Texas: Stata Corp., 1998.
- Sudman, S., N. Bradburn, and S. Schwarz. *Thinking about Answers*. San Francisco, CA: Jossey-Bass, 1996.
- Swartz, K. and J. Purcell. "Letter: Counting Uninsured Americans." *Health Affairs*, Vol. 8, No. 4, pp. 193-197, 1989.
- Swartz, Katherine and Patrick J. Purcell. "Letter: Counting Uninsured American." *Health Affairs*, Winter 1989, pp. 193-196.
- Swartz, Katherine. "Changes in the 1995 Current Population Survey and Estimates of Health Insurance Coverage." *Inquiry*, Vol. 34, No. 1, Spring 1997, pp. 70-79.
- Swartz, Katherine. "Interpreting the Estimates from Four National Surveys of the Number of People without Health Insurance." *Journal of Economic and Social Measurement*, Vol. 14, 1986, pp. 233-243.
- Tourangeau, R., G. Shapiro, et al. "Who Lives Here? Survey Undercoverage and Household Roster Questions." *Journal of Official Statistics*, Vol. 13, No. 1, pp. 1-18, 1997.
- U.S. Bureau of the Census. "1990 U.S. Census." Washington, DC, 1990.
- U.S. Department of Health and Human Services, Health Care Financing Administration. *HCFA 2082 Reports*, various years (1992, 1993, 1994, 1995).
- U.S. General Accounting Office. "Uninsured Children and Immigration, 1995." Publication No. GAO/HEHS-97-126R, Washington, DC: GAO, 1997.
- United States Department of Commerce Bureau of the Census. "Health Insurance Historical Tables." [Available at <http://www.census.gov/hhes/hlthins/historic/> (updated December 2000).]

United States Department of Commerce Bureau of the Census. "Survey of Income and Program Participation (SIPP) Quality Profile." [Available at <http://www.census.gov/sipp/> (undated).]