

## 1.0 Introduction

To meet the demand for information relevant to current health care policy issues, the Federal government has sponsored two national household sample surveys of the utilization of health services and the related expenditures for the care received. The first of these surveys, the National Medical Care Expenditure Survey (NMCES) sponsored by the National Center for Health Services Research (NCHSR), with support from the National Center for Health Statistics (NCHS), covered health care utilization and expenditures by the U.S. non-institutional population during calendar year 1977. The second survey, designated the National Medical Care Utilization and Expenditure Survey (NMCUES), is currently underway and is intended to provide similar comprehensive data for the U.S. non-institutional population for 1980 and for Medicaid eligible families in four states, California, Michigan, New York and Texas. NMCUES is jointly sponsored by NCHS and the Health Care Financing Administration (HCFA).

Both NMCES and NMCUES are panel surveys in the sense that the data are collected for the year of interest by a series of periodic interviews with the initial sample of households. It is important to recognize that the principal purpose of the repeated interviews at appropriate intervals for these surveys is to improve the quality of the data. In both surveys, the key data items include the details of each dental, doctor, clinic or emergency room visit and each hospital stay, including dates and services received; the charges for the health care services received; prescribed medicines purchased and their costs; other medical expenses and finally the source of payment for the care received, that is how much was paid out-of-pocket by the family, and how much was paid by an insurer or other third party, whether public or private.

The purpose of this paper is to discuss several of the methodological issues which arise in the design of medical care expenditure surveys of families and individuals under the assumption that the reference period of interest is of sufficient length to require a panel approach. The panel designs of the NMCES and NMCUES provide an excellent focus for the discussion. The issues include choice of observation unit, coverage problems, frequency and mode of interview, techniques for reducing reporting errors, record check surveys and analysis issues.

## 2.0 Observational Units

### 2.1 Definition

Health care researchers are concerned with the use of health services and the associated expenditures at both the family and the individual level. Interest in family level data arises because decisions to seek or use particular health care services and mode of payment for the care received are often family decisions and family responsibilities. The definition of a "family," therefore, has special importance for medical care expenditure surveys. In both the NMCES and NMCUES panel surveys the family is the primary observational or reporting unit (RU).

Instances in which a family member is temporarily away adds complication. The principal example occurs for sample families with an unmarried child away at school. NMCES and NMCUES include such children less than 22 years of age as family members and require that they be interviewed and their data linked to their family data. Special data collection problems arise in panel surveys which include students living away from home since they may change their place of residence several times a year.

The major complicating factor for panel surveys of families, however, is due to the changing structure of families. Family changes occur continuously throughout the reference period due to birth, death, marriage, divorce or other separation. They often result in the creation of new families, accompanied by changes in address. The changes are often not known until the next interview is attempted with a sample RU. The interview process in a panel survey of families must permit review and comparison of the family structure as previously reported and identification of new RU's. Further, procedures for assignment and implementation of interviews with the new RU's are necessary.

Additions to sample families may occur by relatives moving in, as well as by births. If such relatives were eligible to be selected into the sample initially, then their data should not be included in any individual level analyses of the initial sample cohort. On the other hand, their data needs to be included in any family analyses. The NMCES and NMCUES surveys designate the initial sample individuals as "key" persons. Persons born during the reference period are also "key" persons. Individuals for whom data were collected in subsequent interviews in order to reflect correctly the family structure of the RU, but who were eligible to be selected into the sample initially, are designated as "non-key" persons. Their data are included in family level analyses, but not in any individual level analyses.

### 2.2 Coverage

The design of a single measurement or one-time survey to provide family and individual health care data is generally straightforward compared to a repeated measurement panel survey. The former is usually confined to those families and individuals in existence on the date of the interview and hence ignores the health care experience of those families and individuals who existed at some time during the survey reference period but did not exist on the date of interview. Populations, whether of individuals or families, are dynamic in several dimensions and they are mobile. Panel surveys of medical care expenditures must be prepared to deal with these aspects, or otherwise suffer from significant incompleteness of data due to inadequate coverage of eligible observation units.

Several kinds of coverage problems arise. In a panel survey designed to provide individual and family data for a given time period, it is important to collect data for all

those individuals and families that are a part of the population of interest at some time during that time period. Special attention must be given, therefore, to the collection of data for those eligible sample persons who were alive on the beginning date of the reference period but who died or otherwise left the population of interest before the initial interview; those sample persons who die or otherwise leave the population of interest after the initial interview, but before the end of the reference period; those persons born into a sample family during the reference period; those persons not a part of the population of interest on the beginning date of the reference period who subsequently join the population of interest during the reference period. It should be noted that persons in their last year of life generally use vastly disproportionate amounts of medical services relative to persons remaining alive.

The NMCES and NMCUES surveys are confined to the civilian non-institutional population of the U.S. Provision was made in both surveys to collect data on eligible sample persons who were born or died during the reference period or who were institutionalized, in the army or out of the country (non-temporarily) on the initial reference date but subsequently joined a sample family. Both surveys attempt to locate all sample RU's/individuals that had moved since the previous interview date and to interview them at their new addresses.

### 2.3 Automated Survey Control System

The general size and complexity of medical care expenditure panel surveys, particularly those aspects concerned with maintaining coverage, with structural changes in the sample RU's and with locating and interviewing RU's that move, requires computer supported monitoring of the data collection and data processing operations. Both NMCES and NMCUES make use of an automated survey Control System to maintain a record of the current status of each individual and RU in the survey; to monitor the flow of data, and to schedule data collection as well as to retain information of an historical nature. To accomplish these tasks, event codes for each survey status or activity, such as "completed interview for current round" or "at data entry," are maintained in each individual's record and updated as the survey instruments progress through the various monitoring points including data receipt, manual edit and coding, document control, data entry, machine edit, and reassignment for the next round of interviewing. The Control System maintains appropriate linkages in order to monitor the response from each participant even if the participant moves to a different geographic location, to a different RU, and/or to a different RU and then back into the original RU. The entire set of time-sequenced transactions which occur to an individual in each interview round throughout the survey are retained in a history file.

### 3.0 Frequency of Interview

#### 3.1 Response Errors

Medical care expenditure surveys place a heavy burden on respondents to recall and report each and every medical care event, and the details of those events with respect to dates, services received and costs, for themselves

and for the other members of their families. Incomplete reporting is clearly a problem with the quality of the data somewhat poorer a) for clinic outpatient visits than for hospital stays or physician office visit (cf. Yaffe and Shapiro) and b) for the elderly and the poor than for those younger and better off economically (cf. Andersen *et al.*). Completeness of reporting medical care events also depends to a significant extent on the length of the recall period. For this reason, NMCES and NMCUES use five interviewing rounds to collect data on medical care utilization and expenditures for a calendar year. The average recall period is approximately three months. Since there is a tendency to recall and report events as having occurred later in time than they actually occurred (i.e., to telescope events forward in time), the use of repeated interviews with specific reference dates provides bounded and hence more accurate longitudinal data. In order to minimize recall errors, NMCES and NMCUES also make use of a calendar/diary on which the respondents are asked to record the dates of all medical care visits by members of the RU and to record all other medical expenses such as purchases of prescription drugs.

#### 3.2 Reporting Unit Summaries

Quite frequently respondents may not know at the time of the interview all the information required in a medical care expenditure survey about each and every medical provider visit. For example, the bill for recent visits to a doctor or outpatient clinic may not have yet been received on the interview date. Or the bill may have been received and the total charge known, but the actual family out-of-pocket cost may not be known because the proportion to be absorbed by health insurance or other third-party payer is still unknown. Subsequent interviews offer an opportunity to obtain information missed in the previous interview(s).

To assist the NMCES and NMCUES respondents and interviewers to provide complete and accurate data on dates of medical care visits, the sources of care, the services received, the charges for those services, the sources of payment and the share of the total charge paid by each source of payment, a computer-generated Summary<sup>1/</sup> of these data items is prepared and mailed to the respondents prior to each interview round. The interviewers also receive a copy and are instructed to review the Summary with the respondent in order to fill in data which were not known during the previous interview(s) and to otherwise update, as necessary, each line of data in the printout for each person in the Reporting Unit. In an effort to minimize reporting errors due to the use of proxy respondents, all adult participants in the RU are asked to review the Summary for missing or incorrect data in their own records prior to the scheduled interview.

The use of a Summary such as described here requires immediate processing of all data collected in each interview round in order to produce an updated Summary for the next interview round. This is not an easy task in large and

<sup>1/</sup> See Appendix A for an example of a Summary.

complex panel surveys such as NMCES and NMCUES. The requirement that a Summary be generated and distributed prior to each interview may, in fact, lengthen the desired interval between interviews by several weeks. This happened in NMCES where interview rounds were initially specified at nine-week intervals. Interviewing, data receipt, data processing, Summary production, and reassignment to the field required a thirteen-week turnaround. Clearly, there is a trade-off vis-a-vis data quality between interviewing each RU more frequently, but not using computer-generated Summaries, and longer intervals between interviews with Summaries distributed prior to each interview.

The choice of the most appropriate interval between interviews in medical care expenditure panel surveys, with or without Summaries, is one of the most important unanswered methodological issues for such surveys. The NMCES data offer a special opportunity to address this issue, since they include information on medical care received by survey participants which was collected from the patient records of the physicians, hospitals, clinics and other health care facilities that provided that care. Thus, the relationship of the proportion of verified data to frequency of interview and the corresponding levels of total survey error can be studied in some depth.

#### 4.0 Mode of Interview

##### 4.1 Personal

Medical care expenditure panel surveys are expensive. Consequently there is pressure to use data collection procedures which are less costly, such as conducting interviews by telephone. However, personal or face-to-face interviews are viewed as necessary for at least some of the interview rounds. Personal interviews are essential in the initial interview to introduce the survey to the sample RU's, to motivate their participation in succeeding rounds and to explain and leave the calendar/diary. Personal interviews in medical care expenditure surveys are also essential for introduction to and in-depth review of the Summary. Consequently, NMCES and NMCUES require personal interviews in the second and fifth (last) rounds as well as the first.

##### 4.2 Telephone

Despite the length and complexity of the core questionnaire, telephone interviews are used in the third and fourth rounds of NMCES and NMCUES. The telephone interviews do, however, proceed somewhat faster than the personal interviews. Although a review of the Summary is included in the telephone rounds, it is confined, for the most part, to updating items not known in previous interview rounds. Line-by-line in-depth review of the Summary requires a face-to-face interview.

There is some evidence that respondents underreport utilization of health care services to a greater extent in telephone interviews (cf. Paul Moore) compared to personal interviews with the same average length of recall period. Since Summary production is costly and also lengthens the interval between interviews, an interesting trade-off would drop the Summary for the telephone interview rounds, decrease the interval between those rounds to about eight weeks, and add a third telephone interview. The cumulative

Summary would be reviewed line-by-line during the final personal interview. Such a trade-off must weigh the overall quality of the data for each data collection scheme with the same total expenditure. Essentially it questions the value of the Summary for the telephone interview rounds relative to a shorter recall period.

#### 5.0 Record Check Surveys

##### 5.1 Medical Provider Surveys

Despite the use of repeated, bounded interviews every 13 weeks, calendar/diaries and Summaries of previously reported health care events, the quality of the NMCES respondent reported data on utilization, diagnoses and expenditures was considered sufficiently suspect as to require a record check survey of those physicians, clinics, hospitals and other medical providers who had supplied health care to a subsample of the NMCES participants. The Medical Provider Survey (MPS) was designed to fill in gaps in the family and individual respondent reported data, particularly with respect to unreported visits, diagnoses, total charges and sources of payment. For example, respondents who were enrolled in public programs such as Medicare and Medicaid frequently do not have access to provider bills or other documents that contain cost and source of payment data associated with their visits. In addition, the MPS data, provide an opportunity to adjust for inaccuracies in the respondent reported data.

Briefly, the NMCES Medical Provider Survey was confined to a sample of those survey participants who had signed Permission Forms for the physicians, clinics and hospitals that had provided care to them during 1977 or were designated as the participant's usual source of care. The Permission Form authorized the medical providers to report to the NMCES project information in their files on the specific participant's 1977 medical care visits, diagnoses, services received, charges and sources of payment. The eligible survey participants were stratified according to diagnostic category and type of medical provider prior to selecting the sample for the MPS.

Under the assumption that the medical provider data, has greater accuracy than the household survey data, various alternative data collection methods for medical care expenditure panel surveys can be compared using the NMCES data. For example, the most appropriate interval between interviewing rounds in relation to accuracy and cost can be determined. Other significant methodological issues which can be addressed include (i) comparison of telephone versus personal interviews, (ii) the marginal increase in overall accuracy achieved by the use of Summaries, (iii) comparison of the quality of the data collected from RU's which report using the calendar/diary versus RU's which do not, and (iv) comparison of the NMCES and NMCUES survey designs with survey designs which collect medical provider data for only the subgroups identified as reporting the least accurate medical care utilization and expenditure data in household panel surveys.

##### 5.2 Health Insurance Verification Surveys

The extent to which medical care costs are covered by private as well as public health insurance is of considerable interest. Health insurance coverage for each participant

is determined in each interview round of the NMCES and NMCUES surveys for each survey participant. The accuracy with which household survey respondents provide information on their health insurance coverage is not known at present. The NMCES and NMCUES include verification of health insurance coverage. These verification surveys should provide some very useful answers concerning the quality of respondent reported health insurance in medical care expenditure panel surveys. The population covered by health insurance is changing constantly as individuals change employers or move in and out of eligibility for coverage under Medicaid. Accurate health insurance coverage data at the family or individual level for a calendar year may be difficult to obtain except in panel surveys with interviews occurring at least as frequently as every three months. The NMCES and NMCUES surveys are hampered to some extent with respect to assessing the net bias in respondent reported health insurance coverage. As in the Medical Provider Survey, in which medical care received from providers not reported by the household survey participants could not be verified, so also health insurance plans not reported in the household survey cannot be verified. The NMCES project attempted to reduce this gap in the data by contacting (with participant permission) the provider named as the usual source of care for those participants not reporting any provider visits during the reference year and by contacting (again with participant permission) the employers for those participants not reporting coverage by any health insurance plan.

## 6.0 Analysis Issues

### 6.1 The Changing Population Base

The principal advantage of the panel survey designs implemented for NMCES and NMCUES over a single or multiple cross sectional surveys is their collection of accurate day-by-day accounts of all health care related events for a national probability sample of families and individuals, thereby providing the ideal setting for the consideration of continuous time models of the health care utilization and expenditure process over a twelve-month interval. Along with these advantages comes attendant complications in the definition of descriptive population parameters in light of the changing population base over time. The definition of simple quarterly and annual health care utilization and expenditure rates per person and per family require one to specify for each quarter and for the year a single population size for the denominator of these parameters. Since the eligible population for these surveys is the civilian noninstitutionalized household population residing in the United States at any time during the reference year, it is clear that the size of this eligible person universe will change daily due to births, deaths, and changing institutional and military status. With the addition of separations, divorces and remarriages, the universe of eligible families is subject to even greater fluctuations over time than the person universe. The most common response to this complication has been to use the population size at the beginning or the end of the reference period as the base for defining rates.

An appealing alternative that is facilitated by the NMCES/NMCUES survey's continuous time monitoring of life events for a probability sample of family reporting units and all the persons that subsequently join these families is to estimate the average number of universe members per day during the reference period. Aside from previously ineligible institutionalized and military personnel who rejoin the population as unrelated individuals with no chance of being interviewed as part of an originally eligible family, all persons who are eligible household population members for any period of time during the reference period have a known probability, say  $\pi(i)$ , of being in the sample. The selection probability for eligible person (i) is the inclusion probability for the housing unit containing person (i)'s initial round family reporting unit or the reporting unit that person (i) subsequently joined. To avoid multiplicity complications, "non-key" persons who join NMCES/NMCUES sample families after the first round of interviewing are not included in the estimation of person based parameters. To estimate the number of eligible persons on day t of the reference period, it is possible, for the NMCES and NMCUES surveys, to define day t eligibility indicators  $E_t(i)$  for all key sample members (i). The one exception to this rule is caused by the failure to collect the reentry date when a previously ineligible person joins a family reporting unit. This oversight should be easy to correct in subsequent surveys. By collecting reentry dates and by rescreening the original sample housing units (HU's) and the associated half-open intervals between the sampled HU's and the next listed HU for the reentry of previous ineligible as unrelated individuals, an unbiased estimate for the number of NMCES/NMCUES eligible persons alive on day t is

$$\hat{N}(t) = \sum_{i \in S} E_t(i) / \pi(i)$$

where the summation above extends over all key persons (i) belonging ( $\in$ ) to the sample (S) of family reporting units at any time during the year. The average of these daily population size estimates over a specified quarter q has the form

$$\hat{N}(q) = \sum_{i \in S} PE_q(i) / \pi(i)$$

where  $PE_q(i)$  is the fraction of days during the quarter q that person (i) was eligible. If  $Y_t(i)$  denotes the number of health care provider visits or the associated expenditures reported for person (i) on day t, then the total number of visits or expenditures experienced by eligible persons is estimated as

$$\hat{Y}(q) = \sum_{i \in S} Y_q(i) / \pi(i)$$

where  $Y_q(i)$  is the aggregate visit count or expenditure total for the days in quarter q that person (i) was eligible. Using

$$WE_q(i) = PE_q(i) / \pi(i)$$

as a sample weight for estimating the average number of eligible persons per day during quarter

q, the estimated utilization/expenditure rate per person can be computed as

$$\hat{R}(q) = \hat{Y}(q) / \hat{N}(q) \\ = \sum_{i \in S} W_{eq}(i) [Y_q(i) / PE_q(i)] / \sum_{i \in S} W_{eq}(i).$$

If U denotes the universe of ever eligible persons (i), then  $\hat{R}(q)$  is a consistent estimator of

$$R(q) = Y(q) / \bar{N}(q) \\ = \sum_{i \in U} PE_q(i) [Y_q(i) / PE_q(i)] / \sum_{i \in U} PE_q(i)$$

which can be interpreted alternatively as having the effect of treating  $PE_q(i)$  as an estimate of the common probability that person (i) will be eligible on any day t during the time interval q. With this alternative interpretation, the inflated variate values  $[Y_q(i) / PE_q(i)]$  would estimate person (i)'s aggregate Y contribution for the entire time interval including ineligible days.

For family level rates, the approach suggested above is complicated by the ability of families to split forming two or more essentially new families. While it is clear that one new family is always created by a split it is not obvious whether the residual portion of the original family should retain its original identity or should be viewed as a new family. Ones ability to associate family characteristics with family level health care experiences would seem to be at the heart of this issue. When family characteristics change enough to influence health care experiences in an important way, then one can argue that the original family should be dissolved and a new residual family identified with the appropriate new family characteristics ascribed. For the NMCES family level analysis, an original sample family will be assumed to exist as long as the family head and spouse remain unchanged. Therefore, when a son or daughter leaves home to set up a separate household (with the exception of unmarried college students less than 22) a new family is born while the original family is assumed to survive. If, on the other hand, the head or spouse leaves the family due to separation, divorce, death, or institutionalization, it will be assumed that the original family died on the date of separation and a new surviving unit was born. Additions of a new family head or spouse due to marriage or the reentry of a previously ineligible person would also result in a head or spouse change and would be considered justification for dissolving the original family and creating a new unit.

A second complication with the family level estimation process for panel surveys like NMCES and NMCUES is the multiple chances of selection for families receiving "non-key" members at some point subsequent to the initial round. By definition, these "non-key" persons were eligible in the initial round and could have been selected at their initial round address leading ultimately to the same NMCES/NMCUES family by a different route. Counting the multiplicity  $m(j)$  for a family (j) in terms of the number of eligible persons in the universe

during the initial round who belong to family (j) at some time during the year and letting  $n_k(j)$  denote the number of initial round sample persons who originally belonged to the initial round family reporting unit (k) that spawned family (j), then providing for the NMCES occurrence where two original sample members from different families married to form a new family, the multiplicity sample weight for estimating the total number of families existing during the year has the form

$$FW(j) = \sum_{k \in S} n_k(j) / m(j) \pi(k)$$

where  $\pi(k)$  is the selection probability for a selected initial round sample family reporting unit that subsequently donates members to family (j). While one can also formulate family weights in terms of family instead of person level multiplicities, we suspect that the person level adjustment factors  $[n_k(j) / m(j)]$  would be more stable and lead to less variable estimators than would family or reporting unit multiplicity adjustments. Following the approach suggested previously for person level rates, the average number of families existing on a given day t would be unbiasedly estimated by

$$\hat{F}(t) = \sum_{j \in S} FW(j) FE_t(j)$$

where  $FE_t(j)$  depicts an existence indicator for family (j) on day t. These existence indicators are clearly a function of the rules established for dissolving initial round families and creating new ones. The head or spouse change rule adapted for NMCES was motivated by the desire to change the family identification when an event significantly altered the families characteristics, but to preserve families when a less serious change occurred so as to maximize the number of families existing for the entire twelve-month reporting period. With this motivation, the average number of families existing during the time interval q with the set of family characteristics identifying a population subgroup c is

$$\hat{F}_c(q) = \sum_{j \in S} FW(j) \Theta_{eq}(j) X_c(j) \\ = \sum_{j \in S} FWE_q(j) X_c(j)$$

where  $\Theta_{eq}(j)$  is the fraction of days during quarter q that family (j) exists and  $X_c(j)$  is a zero-one indicator for family subgroup c. If  $Z_q(j)$  denotes the aggregate number of visits or expenditures for family (j) during quarter q and  $FWE_q(j)$  is the eligible days adjusted weight for family (j) during quarter q, then the associated per family utilization or expenditure rate is the  $FWE_q(j)$  weighted mean of the inflated family totals  $\{Z_q(j) / \Theta_{eq}(j)\}$ . For a rate specific to families of type c, the weighted averaging is simply restricted to type c families.

A second class of descriptive parameters that are affected by the changing population base are distributions of persons and families where the persons or family's position in the distribution depends, for example, on the total number of visits or dollars spent for health care during the interval. In this distributional case the solution analogous to that advocated above for rates somehow seems less appealing.

The analogous solution would suggest classifying families, for example, on the basis of the eligible day inflated aggregates

$$Z_q^*(j) = Z_q(j)/\theta E_q(j).$$

If  $\hat{G}_q(\vec{z})$  denotes a vector of  $r$  zero-one indicators for  $r$  intervals on the quarterly visit count total or expenditure total scale, then the associated distribution estimator would be.

$$\hat{G}_c(q) = \sum_{j \in S} FWE_q(j) G_q[Z_q^*(j)] X_c(j) / \sum_{j \in S} FWE_q(j) X_c(j).$$

Extending the concept above to estimating annual distributions, one is lead to the distasteful prospect of inflating data to yearly levels and imputing a corresponding position in the distribution based on only one or two months of real data. An alternative approach that seems less damaging might be to use the  $FWE_q$  weights with an inflated  $Z_q$  value obtained by a hot deck direct imputation procedure. At the persons level, this process could be enhanced by collecting data in the last round for people who die or are institutionalized after the end of the year but before the final round interview date. Persons who die or become institutionalized during the twelve month reference period could then have their records artificially completed by adding data for the missing months from a selected donor among the final round deaths and institutionalizations. This process could be viewed as approximating the last twelve months of health care experience prior to becoming ineligible for each person who becomes ineligible during the reference period. If, for example, a person dies in March one would ideally select a similar person who died after December but prior to the round 5 interview and impute their experience to the original decedent with the idea that one is approximating the last twelve months of health care experiences for those who do not survive the entire time interval as eligibles. A similar strategy could be applied to data gaps resulting from births by selecting similar donors who were born at roughly the same time the previous year and using their data from the early months of the survey period to extend the newborn's data record to twelve months. A similar imputation combined with the  $FWE_q$  weights could be used to produce artificial twelve-month data histories and a means of weighting them together to properly reflect parameters of the average population existing during the reference period. Short of some attempt to make all survey respondents have the same twelve-month period of risk for health care episodes, one must segregate out those persons and families with less than twelve months eligibility for separate analysis. While the separate analysis approach is not unreasonable for descriptive purposes, the imputation strategy would seem inherently superior for stochastic modeling of the health care utilization and expenditure process. Since the preferred approach elaborated in the following section for dealing with data gaps resulting from survey attrition will be to impute the missing data, the extra chore of imputing for the comparatively small set of eligibility related gaps does not seem overly ambitious.

## 6.2 Survey Attrition

The NMCES survey experienced 12.5% attrition in the sample of individuals between the first round of data collection and the fifth round where interviews were conducted for the reference period ending on December 31, 1977. To adjust for the potential bias associated with different patterns of attrition across 32 age (8) by race (2) by sex (2) subgroups, a two stage weighting class and post-stratification adjustment was made to the survey weights. The weighting class phase of this adjustment had the effect of calculating an initial quarterly weight for a partially responding person (i) of the form

$$WR_q(i) = PR_q(i)/\pi_A(i)$$

where  $\pi_A(i)$  was an adjusted version of the inclusion probability for the housing unit containing person i. The adjustment to  $\pi(i)$  amounted to multiplying the initial housing unit inclusion probability by the initial round housing unit response rate in the area segment containing person (i). The  $PR_q(i)$  quantity represents the fraction of days during quarter q that person (i) responded. When these initially adjusted person weights are summed within one of the 32 age race sex cells (c), an estimate of the average number of persons of type c responding per day results; namely,

$$\hat{NR}_c(q) = \sum_{i \in S} WR_q(i) X_c(i).$$

This average number of responding persons per day is divided into the corresponding average number of eligible persons per day derived from the original eligible days adjusted weights

$$WE_q(i) = PE_q(i)/\pi_A(i).$$

If  $\hat{NE}_c(q)$  denotes this average number of eligible persons per day in subgroup c during quarter q, then the combined ratio estimator of the average daily response rate

$$\hat{\rho}_c(q) = \hat{NR}_c(q)/\hat{NE}_c(q)$$

was used as a further adjustment assuring that the weights

$$WR_q^{(1)}(i) = PR_q(i)/\pi_A(i)\hat{\rho}_c(q)$$

will sum to the average number of eligible persons per day  $[\hat{NE}_c(q)]$  within each of the 32 age by race by sex cells. As a final step these  $WR_q^{(1)}(i)$  weights were ratioed up to corresponding 1977 quarterly CPS population counts for each cell c based on the middle month of the quarter. If  $CPS_c(q)$  represents these current Population survey counts, then the final post-stratification adjustment factor  $g_c(q) = \hat{NE}_c(q)/CPS_c(q)$  was applied to the weight  $WR_q^{(1)}(i)$  to form

$$WR_q^{(2)}(i) = PR_q(i)/\pi_A(i)\hat{\rho}_c(q)g_c(q).$$

While the While the weighting class and post-stratification adjustments could have been combined into a single factor of the form

$$\hat{\rho}_c(q)g_c(q) = \hat{NR}_c(q)/CPS_c(q)$$

there was interest in the relative size of the two factors, with the second suggesting the NMCES under coverage rates relative to CPS counts for the 32 cells. The smallest of these NMCES/CPS coverage rates was only slightly over .5 for black males in the less than 25 age bracket.

The implication of the weighting strategy described above for utilization and expenditure rates is to use a weighted combination across cells  $c$  of the total number of respondent reported visits or expenditures divided by the average number of responding persons per day in cell  $c$ ; that is

$$R_{(2)}(q) = \sum_{c=1}^{32} \text{CPS}_c(q) \hat{RR}_c(q) / \sum_{c=1}^{32} \text{CPS}_c(q)$$

where the cell specific rates for respondents have the form

$$\hat{RR}_c(q) = \hat{YR}_c(q) / \hat{NR}_c(q)$$

with the numerator estimating the total number of visits or expenditures that would be reported by partial respondents for cell  $c$  if a census was conducted. In terms of the original responding day adjusted weights  $WR_q(i)$ ,  $RR_c(q)$  has the form

$$\hat{RR}_c(q) = \frac{\sum_{i \in S} WR_q(i) X_c(i) [YR_q(i) / PR_q(i)]}{\sum_{i \in S} WR_q(i) X_c(i)}$$

with the value in square brackets reflecting an inflation of the responding day contributions  $YR_q(i)$  by the reciprocal of the fraction of responding days in the quarter. As with the inflation based on the fraction of eligible days during the quarter alluded to in the previous section, this responding days adjustment while reasonable for rates is not very appealing for estimating distributions since by analogy person  $(i)$  would be categorized for placement in the associated annual histogram on the basis of the simple inflation no matter how small the responding day fraction  $PR(i)$ .

A direct data imputation strategy would seem to be a superior approach for dealing with partial data gaps due to survey attrition than the responding day weight adjustments displayed above. A carefully executed hot-deck type imputation would fill in the missing time periods with actual visit and expenditure records borrowed from a randomly selected donor drawn from a poststratum of individuals who completely respond for their entire period of eligibility and are similar in individual characteristics and initial health care experiences to the persons in need of data imputations. For example, one could use at a minimum the 32 age by race by sex cells mentioned previously as explicit poststrata and in addition sort the donors and recipients within poststrata on the basis of total visits and expenditures for comparable time periods. These selections of donors could proceed in waves beginning with imputations for second quarter dropouts, with the sorting of donors and recipients within poststrata based on first quarter visit counts and expenditure totals. Initial round respondents interviewed in February and early March who drop out in the second round might reasonably be excluded with compensating weight adjustments for their total nonresponse so as to avoid imputations for first quarter dropouts.

A new randomized hot-deck selection algorithm has been recently adapted at RTI by Dr. Brenda Cox from a sequential PPS minimum replacement selection routine developed by Dr. James Chromy. This weighted sequential hot-deck algorithm determines the number of times a potential donor is to be used at random in such a fashion that over repeated donor selections the weighted mean for the poststratum based on both the actual and imputed responses is equal to the weighted mean for complete respondents. The association of selected donors with recipients is made in the order of the sort imposed on both the donor and recipient lists within poststrata thus preserving some of the implicit poststratification inherent in the standard nearest neighbor sequential hot-deck. While one might occasionally be able to demonstrate that the sorting variable was so highly correlated with the missing data that a nearest neighbor imputation would minimize imputation bias, it is felt that in the great majority of applications, nearest neighbor imputation will not be as important as recognizing the unequal selection probabilities of the donors and recipients so that over repeated imputations, the weighted distribution of imputed and actual values would estimate unbiasedly the population distribution for all potential donors in the poststratum. The randomized selection of donors minimizes the number of times a respondent is selected to provide imputations consistent with achieving the desired expectation over repeated imputations. The randomized nature of the algorithm is also conducive to making multiple imputations within specified subsample replicates allowing one to partition the variance of imputation based statistics into between sample and between imputation within sample components. A paper detailing Cox's weighted sequential hot-deck is being presented at the 1980 ASA meetings. The application of this hot-deck algorithm and a hybridized combination of it and a regression imputation method are discussed in the context of medical provider record checks in the following section.

### 6.3 Provider Record Check Data

It is generally conceded that medical provider record check surveys that use the household respondents to identify the frame of providers eligible for record checking can improve the quality of household reported utilization and expenditure data. This assumes, of course, that the medical provider data can be easily and accurately matched to household respondent data. Some improvement seems possible over the NMCES experience, particularly since the MPS data were collected after the household survey data collection was completed. More work is needed to develop effective survey methods designed to get the providers involved early in the survey to assist in the matching of household reported and provider recorded visits. Matching is essential to the proper identification of unreported and overreported visits. If discrepancies between household and provider reports are identified prospectively during a survey, greater opportunity for needed reconciliation and adjudication exists. Permission to contact a provider therefore needs to be secured as soon as a summary containing the reported visit is available.

Record checks based on subsamples of household respondents are both feasible and potentially optimal since data collection and processing cost reductions generally trade-off favorably with sampling induced variance increases. As indicated previously, the NMCES medical provider survey sample was selected at a 100 percent rate from certain infrequent diagnosis categories and from outpatient clinic and emergency room only visitors. The designed subsampling rate was 40 percent for hospital inpatient and other outpatient facility users and 20 percent for physicians office visitors only. Subsequent MPS survey optimization studies conducted by the second author suggest that aged hospital users and nonaged poor outpatient facility users should also be oversampled heavily.

With an overall MPS subsampling rate of less than 50 percent of the household respondents, the issue of how to supply provider based utilization, expenditure, and diagnosis data to the analysts was resolved in favor of imputing provider data to the full household survey file of visit level responses. In this double sampling application of imputation methods the donor sample of provider responses is a valid probability subsample of the full household sample. With this in mind, one can show that an MPS subsample based regression equation that predicts provider reported expenditure amounts per visit as a function of household reported expenditures can be used to impute provider expenditures to non-MPS household responders in such a manner that the resulting full sample estimates of expenditure means will be equivalent to a double sampling regression estimator. This is accomplished by using the household sample weights in the regression model and fitting within each MPS poststratum so that the sum over the MPS subsample of residuals from the predicted expenditures is zero. This result allows one to recast the double sampling regression estimator for a total which amounts to the full sample weighted sum of predicted values as the sum of observed provider values for MPS subsample members and the sum of predicted values for non-MPS household survey members.

While these regression based expenditure predictions can be used to form consistent double sampling estimators of per person or per visit expenditure rates, they should not be used to estimate distributions since the latter would be biased, possibly substantially, by the shrinkage in variability induced by elimination of the natural deviations from regression. To remedy this situation a residual value could be imputed in addition to a predicted mean by selecting donors from the MPS sample with similar characteristics and similar predicted expenditures. In this application of Dr. Cox's weighted sequential hot-deck procedure, the model based predicted values for the total charge and various disaggregates should be ideal sorting variables. With the randomization inherent in Dr. Cox's algorithm, the weighted distribution of persons by expenditure level computed with imputed and actual provider measurements is the desired MPS subsample based distribution estimator. Using the regression equations to predict hot-deck cell means with unimportant interaction effects eliminated from the model should provide more precise imputation based statistics than would a

direct application of hot-decking provider expenditure amounts. The direct hot-decking procedure is being used to impute MPS data for non-MPS sample recipients since fitting the regressions is not required.

The imputation of health care visits reported by the NMCES medical providers that were not matched to household reported visits is also being accomplished using Cox's weighted sequential hot-deck method. After poststratifying the MPS sample donors and non-MPS sample recipients on the basis of relevant demographic variables and the types of health care utilization reported in the household survey, MPS sample donors will be selected within poststrata from lists sorted by numbers of household reported visits. When a selected MPS donor has provider reported visits that do not match with a household reported visit, such provider reported visit records will be imputed to the non-MPS sample recipient.

The third NMCES Medical Provider Survey related imputation task that is currently underway involves the assignment of medical provider based diagnoses for household reported visits not covered by the MPS. The weighted hot-deck algorithm will be used in this instance to select matched household and provider visits from poststrata formed in terms of relevant demographic variables and 60 to 65 major groups of household reported diagnoses. Within these demographic by major diagnostic category poststrata, the matched visits will be sorted according to more detailed diagnostic codes derived from the household responses. With the selection of a matched visit donor, the associated provider reported diagnostic code will be imputed to the non-MPS sample recipient visit. While the decision to impute MPS data to the full NMCES household sample was largely predicated on the resulting convenience for data analysts, a valuable by product of this strategy will be the ability to make corresponding cold-deck type imputations of 1977 NMCES medical provider data to the 1980 NMCUES household visits. No medical provider record check was planned for the NMCUES survey on the assumption that the associations between household and provider reported data will not have changed substantially between 1977 and 1980. The two step regression/hot-deck residual procedure should reflect trends in expenditure levels better than the direct hot-deck procedure. Adjustments for trends in expenditure levels are possible, of course, with the direct hot-deck procedure.

#### 6.4 Panel Data Analysis Modes

In addition to the descriptive analyses aimed at utilization and expenditure rates and distributions alluded to in previous sections, panel survey designs like NMCES and NMCUES provide for the modeling of short term trends in utilization and expenditure rates. The longitudinal measurements for families and individuals facilitate the use of multivariate growth curve analyses recognizing that the monthly or quarterly experiences of sample families are correlated from one period to the next. The methods that have been developed recently for fitting linear models in the context of complex probability sample designs can be applied directly to fit polynomial growth models to utilization rates per person and average expenditures per visit



over monthly or quarterly periods. Another class of important models that can only be explored with panel data such as that produced by NMCES and NMCUES are the continuous time markov and semi-markov renewal processes that show considerable promise for predicting the dynamics of health care utilization patterns over time.

REFERENCES

Andersen, R., Kasper, J., Frankel, M. R. and Associates. Total Survey Error. San Francisco, California. Jossey-Bass, Inc. 1979.

Yaffe, R. and Shapiro, S., "Reporting Accuracy

of Health Care Utilization and Expenditures in a Household Survey as Compared with Provider Records and Insurance Claims Records." Paper presented at the Spring Meeting of the Biometric Society, Eastern North American Region, New Orleans, LA, April 1979.

Moore, Paul R., "Comparisons of Personal and Telephone Interview Data from the National Medical Care Expenditure Survey. Task 2: Comparison of Telephone and Personal Interviews," prepared for National Center for Health Statistics under Contract No. 233-78-2102, April 1979.

Appendix A

NATIONAL MEDICAL CARE UTILIZATION AND EXPENDITURE SURVEY  
SUMMARY OF RESPONSES - ROUND 1

HEALTH CARE SERVICES FOR JOHN SMITH FOR THE PERIOD 01/01/80 - 03/12/80

RU ID # 7654321  
PID 1234567

①	DATE OF CARE	TYPE OF VISIT OR SERVICE	MEDICAL PERSON OR PLACE AND ADDRESS	SERVICES RECEIVED	--CHARGE INFORMATION--	
					SOURCE OF PAYMENT	AMOUNT
①	01/05/80	DENTAL VISIT	DR. SAMUEL JONES	FILLINGS (02) FLUORIDE TREATMENT ----- -----	FAMILY ----- -----	\$35.00 ----- -----
					TOTAL CHARGE	\$35.00
②	02/18/80	HOSPITAL OUTPATIENT	EAR, NOSE, & THROAT CLINIC WAKE MEDICAL CENTER RALEIGH, NC	DIAGNOSIS/TREATMENT X-RAYS LABORATORY TESTS ----- -----	FAMILY BC/BS OF NC ----- -----	20% 80% ----- -----
					TOTAL CHARGE	NOT KNOWN
③	02/21/80	MEDICAL VISIT	DR. JANE GREENE RALEIGH, NC	GENERAL CHECK-UP ----- -----	FAMILY ----- -----	\$45.00 ----- -----
					TOTAL CHARGE	\$45.00
④	02/18/80	PRESCRIPTION	AMPICILLIN	1 TIME	FAMILY BC/BS OF NC ----- -----	\$6.82 NOT KNOWN ----- -----
					TOTAL CHARGE	\$6.82
⑤	03/12/80	HEALTH INSURANCE	BC/BS OF NC ----- ----- -----	Q7 - PRIVATE PLAN		