Nancy A. Mathiowetz, NCHSR
5600 Fishers Lane  Rockville, MD  20857

Survey research questions are often interested in quantitative autobiographical information that is difficult for a respondent to recall. Although there is a large body of literature which has documented the existence of response errors and correlates of the error, the discussion of theory related to response error has been limited. More recently, methodologists have turned to the theories of cognitive psychology to address questions of recall, looking at factors as far-ranging as the effects of cueing mechanisms on the retrieval of memories and the effect of present mental state on retrospective recall. The present research will review a small piece of the cognitive literature and test the applicability of these theories to a specific type of response error: errors of omission.

## Factors that Affect the Recall of Episodic Memories

Experimental work completed in the field of cognitive psychology coupled with the research related to non-sampling errors in survey research have generated a list of factors which affect the quality of  retrospective recall. These factors include:

1. Events and information encoded since retention, including the introduction of new and possibly conflicting information;

2. Length of time between the occurrence of an event and the recall of that event;

3. The salience of the event; and

4. Demographic characteristics of the respondent such as education and age.

Intervening Events. Classical interference and information processing theories suggest that as the number of similar events occurring to an individual increases, the probability of recalling any one of these events declines. In psychology experiments, the serial position of a word in a list affects the probability of recalling the word. Words that appear early (primacy effect) and late (recency effect) are more likely to be recalled than words in the middle.

Time and Length of Recall. Intuitively, the most obvious aspect of forgetting is that we recall more and more poorly with the passage of time. McGeoch (1932) was the first to suggest that the hypothesis that memories fade because of the passage of time was theoretically sterile and wrong. He showed that with the amount of time since learning held constant, it was possible to experimentally vary the amount of forgetting by manipulating what went on during the intervening time. His work forced researchers to reevaluated the existence of a correlation between time and memory as an indication that time is the causal factor in memory loss.

From the previous discussion concerning the effect of intervening events, it is easy to suggest that it is not the passage of time that affects memory, but the fact that the longer the period of time that exists between the occurrence of an event and the retrieval of information about that event, the more likely that a number of similar events or the encoding of conflicting information has occurred.

Salience of an Event. Salience as defined by Webster is a characteristic indicating prominence or emphasis. The salience or importance of an event is hypothesized to affect the strength of a memory trace. The more salient the event, the stronger the memory trace and subsequently, the less effort or search of memory necessary to retrieve the information. This should make salient events subject to lower rates of errors of omission. However, this characteristic of salient events makes them more subject to two types of response error, forward telescoping and overestimation.

Demographic Characteristics. In both studies of response error in survey reports and psychological experiments, education and age have been correlated with errors (e.g. Loftus, 1979; Cannell, et al, 1977). The cognitive functions behind these correlations are not well understood. For example, do older individuals forget new information more readily than younger individuals and if so, is the problem related to encoding new information or retrieval of available information? Some evidence suggests that acquisition but not retention is affected by age (Tallard, 1968) while other research indicates that age adversely affects the retrieval process (Craik, 1968).

Education has been shown to be related to high levels of performance in free-recall tasks (e.g. Tulving, 1968). Performance appears to be related to the individual's ability to organize retrieval cues and list items into highly efficient associative structures. The observed education/recall correlation has been hypothesized as related to the "problem solving" aspect of higher education. Episodic memory provides the source of analogies between the current problem and others the individual has solved in the past. The ability to solve problems sharpens an individual's ability to efficiently search episodic memory for parallel events or problems.

## Theories to be Investigated in Present Research

The goal of the present research is to address the applicability of cognitive theories to the problems of response error in surveys. The intent is to determine whether the patterns of response error in surveys parallel some of the findings resulting from laboratory research. Among the questions to be addressed are:

1. Is there evidence for an effect of time (length of recall) after controlling for intervening events? Previous methodological

investigations have supported a length of recall:omission rate correlation, but none of these studies have controlled for intervening events.

2. What role do intervening events have in affecting rates of omissions.

3. Does saliency affect the reporting of a single event or groups of related events?

4. How are cognitive functions affected by age or eduction? Are patterns of response error predictable by age or education?

The reader is cautioned that although the discussion to follow focuses on response error as a function of memory and cognitive processes, several alternative explanations are possible. Unlike most memory experiments in psychology, the task facing survey respondents involves the recall of facts or events which may have characteristics of social desirability. Failure to report an event may reflect misunderstanding of the question, concern with self-presentation in an unfamiliar social setting (the interview), as well as true recall error.

Research Design

The data to be used for testing the applicability of hypotheses generated from cognitive theory are from the State Medicaid Sample of the National Medical Care Utilization and Expenditure Survey (NMCUES) (Bonham, 1983). NMCUES was designed to collect data about the U.S. civilian noninstitutionalized population during calendar year 1980. Information was obtained on health, access to and use of medical services, associated charges and sources of payment, and health insurance coverage. The study, cosponsored by the National Center for Health Statistics and the Health Care Financing Administration, consisted of two data collection components. The national household component consisted of approximately 6000 randomly selected households. The State Medicaid household (SMHS) component consisted of about 4,000 households selected from the Medicaid eligibility files in California, Michigan, New York, and Texas. Administrative records containing information about health care utilization and payments for members of sampled households in the SMHS were obtained from each of the states. The data used in the present research focuses on the SMHS sample and the corresponding administrative records. A more detailed discussion of the sample design can be found in Folsom and Iannacchione (1980).

The first round of data collection was conducted as a personal interview. Interviews were conducted from February through April of 1980, with respondents providing information for the period from January 1, 1980 to the date of the interview. This resulted in a reference period ranging from four to sixteen weeks.[1] Data for all persons related to the selected Medicaid case member were collected. Persons aged 17 and older were encouraged to report for themselves; information for all persons under the age of 17 was collected by proxy. The

household level response rate for the first round of data collection was 86.7 percent.

The field procedure by which data were collected for all members of the household resulted in a mix of self and proxy reporting for adults. Those adults who were present at the time of the interview are by definition self reporters. These self reporters also served as informants about those who were absent from the household at the time of the interview. The analyses presented here will be limited to self reports, to eliminate the problems of confounding response error with knowledge issues.

The questions of interest for this research are the respondent's reports of visits to ambulatory health care providers. Ambulatory health care includes visits made to emergency rooms, hospital outpatient departments as well as visits to office-based medical providers. The term "medical providers" includes persons such as chiropractors, speech therapists, faith healers, psychologists, and nurses, as well as medical and osteopathic doctors.

Administrative Records and Matching Procedures

Claims for the four states in the SMHS sample were abstracted in the fall of 1981, after which approximately 90% or more of all bills for 1980 would have been filed and processed. Although the design of the records system and the information available from each of the states varied, the discussion to follow applies to all four states.

All links between the claims data and the household survey reports were completed manually. Three basic combinations of the variables from the household reports and the claims data were used in the matching process. These combinations were:

1. Type of service and provider.
2. Type of service and date of visit.
3. Provider name and date.

Within blocks defined by person identification numbers, preference for matching was given to those claims and survey records which matched on all three variables--provider, type of service, and date. If only two characteristics matched, priority was given to the combinations in the order listed above. The details as to when a service type or date of service were considered an exact match are provided by Smith (1983).

After the matching work was completed, unmatched claims and unmatched survey reports remained. The assumption concerning unmatched claims is that these claims represent legitimate visits to medical providers that were not reported by the household respondent. Unmatched survey reports were examined with respect to source of payment. If Medicaid was reported as a source of payment for these events, the event was, in most cases, classifed as an overreport by the survey respondent. If Medicaid was not listed as a source of payment, and therefore no Medicaid claim would be expected to match the event, the event was considered a legitimate visit to a medical provider and the report retained in estimates of utilization and expenditures.

The data to be used in this research has several limitations. These limitations include:

1. The data are for respondents who are on Medicaid from a sample of four states, California, Michigan, New York, and Texas. To the extent that these individuals differ in their reporting of visits to physicians, emergency rooms, and outpatient clinics then the general population, the inferences drawn from the analyses may be biased.

2. There is little documentation concerning the matching process and no assessment of the level of matching error for these data. In the absence of this documentation, all errors will be attributed to the respondent, resulting in an overestimation of response error.

3. Documentation of each states' Medicaid claims files raises questions concerning the completeness of these files (Corder, et al., 1984). Incomplete claims files would result in an an underestimation of omissions (reduction in the number of unmatched claims). Completeness of the claims data is dependent upon the amount of time between the date of the event and the abstraction of the claims data. Since the present analysis is limited to the first round of data collection, which was completed approximately eighteen months prior to the abstraction of the claims data, this concern should be less important than if data collected in later rounds were included.

4. The ability to assess error only for Medicaid eligible events leads to a conservative estimate of the actual level of omissions, although most likely will not bias the pattern or correlates of omissions.

These concerns do not suggest that the data are not useful in attempting to understand the relationship between cognitive measures and response error. Rather the limitations suggest that this research should be seen as an exploratory rather than a confirmatory analysis.

## Characteristics of the Data

The data for the NMCUES analyses are limited to self-reporters in the first round of data collection. This restriction results in a sample size of 1280 individuals. The event level data, that is the file where the unit of analysis is a medical event, consists of over 7500 events corresponding to these 1280 individuals. The events are classified according to the source of the data: (1) claim information only; (2) survey information only; and (3) matched interview and claims data. Although the data provide an opportunity to examine both patterns of omissions as well as overreports, the analyses will be limited to omission rates. The restriction results in a database consisting of approximately 1200 individuals with 4700 ambulatory health care visits, consisting of events reported only in the claims data and those matched events reported by both the household and the claims.[2]

Of the 1280 respondents represented in the database, 15.8 percent are men, 84.2 percent are women. The disproportionate percentage of women self-respondents is typical of health surveys which use one respondent to report for the entire household (Mathiowetz and Groves, 1985). The median age of the respondents is 50.9 years; the youngest self-respondent is 17, the oldest is 95. The median education level is less than a high school diploma (9-11 years of education); approximately 4 percent of the respondents had no formal education and 9 percent had at least some college.

## Omission Rates over Time

Figure 1 presents the percent of events reported by the respondent according to the number of weeks between the interview and event date. As evidenced in Figure 1, the rate of decay for reports of health care visits is quite rapid. For events occurring during the week of the interview, the rate of omissions is approximately 2 percent. Events occurring six weeks prior to the interview are subject to an omission rate of almost 30 percent. Beyond week nine or ten, the omission rate is approximately 40 percent.

Three models which attempt to fit the pattern described in Figure 1 are presented in Table 1. The dependent variable for Table 1 (and all tables presented in this paper) is a dichotomous variable in which a value of "1" indicates that the visit was reported by the respondent and a value of "0" indicates that the event was only reported in the claims data. The marginal probabilities therefore correspond to the probability of a visit being reported by the respondent. Logistic regression is used for all of these models and the standard errors presented in the tables have been inflated to reflect the nonindependence of obtaining information for several visits from one respondent.[2]

The three models presented in Table 1 are (1) continuous form of retention period, measured in weeks; (2) the negative exponential transformation of the retention period measure; and (3) a dummy coded recall measure with three segments. The exponential form of the recall period fits the data less well than either the continuous form or the dummy coded measures. Based on the concurrent goal of maximizing fit and maintaining a parsimonious model, the continuous form of the recall period measure will be used in multivariate models presented later in this paper.

## Event Characteristics and Demographics Related to Forgetting

Previous research on the underreporting of health events in household interviews has cited both characteristics of the event and demographic characteristics of the individual as related to rates of omissions (e.g. Cannell, et al, 1977). These studies have found such factors as type of medical condition reported as the reason for the visit, respondent's age and education to be related to the percent of visits not reported. Similar analyses are presented

here in an attempt to replicate earlier findings.

Several aspects of a medical provider visit could be considered as measures of saliency. These might include the cost of the visit, the type of condition related to the visit, the person's overall health status (in that a healthy person may see a visit as more of an anomoly than a person with a poor health status), or the site of medical care. The respondents in the present study are all Medicaid recipients; charges associated with receiving medical care are therefore not relevent. Medical condition associated with the utilization of medical services can serve as a proxy measure of salience. If salience is viewed as a dimension measuring impact and significance of an event, then clearly health visits for life-threatening illnesses would be considered more salient than visits for routine check-ups or a minor irritation such as a sore throat or the flu.

Two of the conditions hypothesized to be related to accurate reporting are malignancies and heart conditions. Although these categories are quite broad, they fit the definition that salient events are emotional and mark a transition point in a person's life.

Place of care, specifically emergency room care versus ambulatory care in all other settings, is also hypothesized to be related to improved reporting. Once again, this hypothesis is based on the belief that emergency room care is likely to represent a more salient event than care received in a physician's office.

Most research on the relationship between demographic characteristics and omission rates has demonstrated that older and less educated respondents are poorer reporters than younger and more educated respondents.

The model presented in Table 2 examines the effects of demographic characteristics, measures of saliency, and retention period simultaneously. The model provides a useful link to earlier models of omission rates. Substantively, the type of condition associated with the medical care visit, site of care, and retention period have the largest effect on whether an event is reported by a respondent. For example, the probability of a respondent reporting an emergency room visit is 32 percent more likely than other types of visits, controlling for all other factors in the model. The marginal probability for retention period (-.03) indicates that for every week the reference period is extended, the rate of omissions increases by 3 percent. The only finding from Table 2 that is somewhat inconsistent with predictive hypotheses is the direction of the age coefficient, although not significant, indicates that older respondents are subject to lower levels of omission rates than younger respondents.

## Effects of Serial Position and Interference

The concurrent effects of retention period, measures of interference and inhibition, salience, and demographics are presented in Table 3. Retention period continues to be strongly related to rates of omissions, even when controlling for other factors hypothesized to affect reporting rates.

Among the measures related to interference and inhibition theories, specifically the total number of health care visits, the spatial relationship among the visits, and an indicator for serial position, only the variable indicating that the event was the last of those occurring during the reference period is related to errors of omission. The last episode of medical care utilization was significantly more likely to be recalled than any other event for the person. Beyond the fact that the last event is the most recent, the serial position is strongly related to reporting, and therefore implicitly, recall ability.

The three measures of the saliency of the event are significant in the model, a finding that implies that the characteristics of an event rather than dimensions of the overall recall task (e.g. the total number of visits) affect recall.

The overall poor fit of the model ($R^2$ = .16) is disturbing. Much of the variance in omission rates remains to be explained. Beyond examining interaction effects, the data provide few additional factors to be used in testing hypotheses related to cognitive theory. Several interaction terms were added to the model (for example, interaction between retention period and characteristics of the event). Although many of these variables had significant coefficients, the small size of the coefficients makes discussion of their substantive importance irrelevant.

The findings from these analyses are somewhat disheartening. The results suggest that with respect to omissions in reports of behavior, most of the variation is unexplained by theory based on cognitive research. The findings confirm those of other studies, most specifically the negative effect of increasing recall period on accurate reporting. However, with the exception of serial position, the indicator that the event was or was not the last event experienced by the individual, the measures related to interference were not related to levels of omissions.

Footnotes

[1] The distribution of persons by month of
interview is as follows: 29 percent in
February, 40 percent in March, 30 percent in
April and 1 percent in May or June.

[2]The standard errors of the regression
coefficients have been inflated to adjust
for the nonindependence of obeservations due
to a single person reporting on one or more
events.  Since the regression models
reported in the paper are logistic
regression models (for which it is difficult
to adjust for design effects), the standard
errors from the logisitic regression models
have been inflated by the ratio of the
variance of the coefficients adjusting for
the design to the variance obtained under
the hypothesis of simple random sampling for
linear regression models.  The largest
effexts were recorded for those variables
that are invariant across events for a
person (e.g. age and education); the highest
ratio was equal to 7.3.

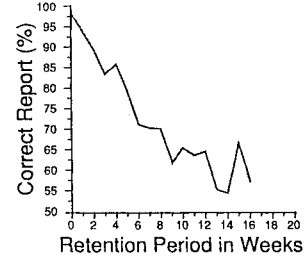Percent of Visits Reported by Retention Period



Table 1.  Logitstic Regression Coefficients:  Effects of Retention Period on Reports of Ambulatory Health
Care Visits.

| | Model I | | Model II | | Model III | |
|---|---|---|---|---|---|---|
| | Coeff. | Marginal Prob.[2] | Coeff. | Marginal Prob.[2] | Coeff. | Marginal Prob.[2] |
| Intercept | 2.29** (.075) | 0.406 | 0.93** (.038) | .165 | 2.11** (.072) | .374 |
| Weeks[3] | -0.16** (.012) | -.028 | | | | |
| EXP(-Weeks)[4] | | | 5.22** (.632) | .925 | | |
| Categorical Weeks[3] | | | | | | |
| <4 | | | | | - | - |
| 4-8 | | | | | -1.12** (.120) | -.198 |
| >8 | | | | | -1.62** (.122) | -.287 |
| Model $\chi^2$ | 344.65[5] | | 273.82[5] | | 344.97[5] | |
| $R^2$ | .07[6] | | .05[6] | | .07[6] | |

171

Table 3. Logistic Regression Coefficients: Effects of Demographic Characteristics, Type of Condition, Place of Care, and Retention Period on Reports of Ambulatory Health Care Visits[1].

| | Coeff. | Marginal Prob.[2] |
|---|---|---|
| Intercept | 1.11** (.185) | .97 |
| Age. | 0.01 (.020) | .002 |
| Education | 0.19 (.175) | .034 |
| Visit Condition | | |
| Malignancy | 1.30 (.718) | .230 |
| Heart Condition | 0.98** (.264) | .174 |
| Health Status: Poor | 0.11 (.546) | .019 |
| Site of Care: Emergency Room | 1.84** (.267) | .326 |
| Retention Period[3] | -0.17** (.011) | -.030 |
| Model $x^2$ | 38.47[3] | 530.55[3] |
| $R^2$ | .01[4] | .10[4] |

[1]Dependent variable is a dichotomous variable indicating whether visit was reported by respondent (1=reported). Numbers in parentheses are standard errors.

[2]Marginal probability calculated as b*(p*(1-p)) where p is the mean proportion on the dependent variable.

[3]Number of weeks between events and date of interview.

[4]Model $x^2$ calculated as difference between model consisting of intercept only and current model.

[5]$R^2$=(model chi-square -2p)/(2L(0)) where p is the number of variables in the model and L(0) is the maxium log-likelihood with only intercepts in the model.

**p < .01

Table 6. Logistic Regression Coefficients: Effects of Retention Period, Interference Measures, Type of Condition, Place of Care, and Demographic Characteristics on Reports of Ambulatory Health Care Visits[1].

| | Coeff. | Marginal Prob[2] |
|---|---|---|
| Intercept | 0.37 (.195) | .066 |
| Retention Period | -0.12** (.011) | -.021 |
| Interference | | |
| Last Event | 2.81** (.331) | .498 |
| Total Visits | 0.00 (.054) | .000 |
| "Cluster" | 0.27 (.215) | .048 |
| Visit Condition | | |
| Malignancy | 1.39* (.722) | .246 |
| Heart | 0.92** (.268) | .163 |
| Health Status: Poor | 0.23 (.560) | .041 |
| Site of Care: Emergency Room | 1.81** (.268) | .321 |
| Age | 0.01 (.020) | .002 |
| Education | 0.22 (.180) | .039 |
| Model $x^2$ | 813.45[3] | |
| $R^2$ | .16[4] | |