

TABLE III.2

## POPULATION, INTERVIEWS AND SCREENING BY SITE

Site	Estimated Population of Included Exchanges	1990 U.S. Census Population	Estimated Population/ Census	Smoker Households Interviewed	Total Screened Households
Baltimore	705,595	736,000	0.96	157	540
Cleveland	596,428	506,000	1.18	141	446
Columbus	867,896	633,000	1.37	167	692
Dallas	1,254,949	1,007,000	1.25	225	980
Denver	532,059	468,000	1.14	107	475
Houston	1,878,734	1,631,000	1.15	317	1,440
Philadelphia	1,552,950	1,586,000	0.98	311	1,209
Portland	540,546	437,000	1.24	107	433
Total	7,929,157	7,004,000	1.13	1,532	6,215

TABLE III.3

## POPULATION, INTERVIEWS AND SCREENING BY STRATUM

Stratum	Estimated Population	Smoker Households Interviewed	Total Screen Households
Baltimore Low-Income	93,722	29	70
Baltimore Remainder	611,873	128	470
Cleveland Low-Income	214,919	57	164
Cleveland Remainder	381,509	84	282
Columbus Low-Income	100,675	22	80
Columbus Remainder	767,221	145	612
Dallas Very Low-Income	39,957	9	21
Dallas Low-Income	175,629	20	104
Dallas Remainder	1,039,363	196	855
Denver Low-Income	72,218	23	67
Denver Remainder	459,841	84	408
Houston Low-Income	357,574	59	229
Houston Remainder	1,521,160	258	1,211
Philadelphia Low-Income	345,996	77	256
Philadelphia Remainder	1,206,954	234	953
Portland Low-Income	33,954	10	35
Portland Remainder	506,592	97	398
Total	7,929,157	1,532	6,215

were supplied by CACI Marketing Systems. It provided census statistics for each tract in each fire-service area (see section IV.B). We used this list of tracts for the five areas where we had no existing tract inclusion list.

On the basis of this comparison, we were able to create a flag record (FSA, for "fire-service area"). The FSA flag has values of 1, 2, or 3, where 1 indicates the record pertains to a household in a tract that is entirely or partly inside the fire-service area,<sup>3</sup> 2 indicates the household was in a tract outside the fire-service area, and 3 indicates that we are unable to tell whether or not the household was inside the fire-service area. Records are marked 3 if we have no tract for the household.

Table III.4 shows the results of an analysis of how often interviewed households were outside the fire-service areas. As expected, more interviews with households outside the fire-service areas occurred in cities for which we had overestimated the population. In such cases, our definition of the fire-service area included telephone exchanges serving a significant number of households outside the area. Overall, 9.8% of the interviewed households were outside the associated fire-service area. Census tracts could not be determined for an additional 4.0 percent of the households interviewed.

Two caveats to the data user are in order. First, the FSA flag indicates where the interviewed household is in a tract that was completely or partially included in a fire-service area. This is not the same as saying that the FSA flag identifies households in the fire-service area, because a household in a tract that is only partly in the fire-service area may or may not be in the fire-service area. These "border tracts" (shown in Table III.5) exist in Columbus, Cleveland and Houston. Second, the analysis resulting in Table III.4 included households for which the tract information was imputed solely on the basis of zip code (see Section VI.A). Data users wanting to analyze data only for households that are known to be inside the fire-service areas should exclude all households in border tracts or with imputed tracts.

---

<sup>3</sup>In all but three sites - Cleveland, Columbus and Houston - all tracts are entirely within the fire service areas. In Cleveland, Columbus and Houston, some tracts are partly inside.

TABLE III.4

## HOUSEHOLDS INSIDE/OUTSIDE FIRE SERVICE AREA (FSA)

Stratum	Household in FSA Tract <sup>a</sup>	Household Not in FSA <sup>b</sup>	Tract Not Determined	Total
<b>Baltimore</b>				
Low-Income	28	0	1	29
Remainder	121	5	2	128
<b>Cleveland</b>				
Low-Income	47	8	2	57
Remainder	62	19	3	84
<b>Columbus</b>				
Low-Income	21	1	0	22
Remainder	115	24	6	145
<b>Dallas</b>				
Very Low-Income	9	0	0	9
Low-Income	20	0	0	20
Remainder	151	35	10	196
<b>Denver</b>				
Low-Income	22	0	1	23
Remainder	61	18	5	84
<b>Houston</b>				
Low-Income	58	1	0	59
Remainder	219	30	9	258
<b>Philadelphia</b>				
Low-Income	72	1	4	77
Remainder	216	4	14	234
<b>Portland</b>				
Low-Income	9	1	0	10
Remainder	90	3	4	97
<b>Total</b>	<b>1,321</b>	<b>150</b>	<b>61</b>	<b>1,532</b>

<sup>a</sup> The tract information was imputed on the basis of zip code for 83 of the 1,321 households in this column.

<sup>b</sup> The tract information was imputed on the basis of zip code for 12 of the 150 households in this column.

TABLE III.5

## TRACTS PARTIALLY INSIDE FIRE SERVICE AREAS

Cleveland	Columbus			Houston			
1051	3.1	74.24	88.22	211	253	446.01	545.12
1052	11.2	74.9	88.25	212	254	448	545.32
1061	19	75.2	92.1	222.01	263	449.2	547.98
1231	25.1	75.31	93.61	222.02	264	450	548.98
1232	26	75.32	93.62	223.01	273	451.32	549
1244	32	75.33	93.71	224.01	322.01	452.01	550
1371	43	75.34	93.73	224.02	334	452.12	552
1413	44	75.4	93.74	224.03	336	452.22	555.01
1922	45	75.5	93.81	226.01	341	529.01	555.12
	51	76	93.86	226.02	354	530.01	556.01
	62.2	77.1	93.9	228.01	359.11	530.02	559.01
	62.3	77.21	94.1	228.02	361	530.03	701.03
	63.1	77.22	94.2	229	362	531.01	701.14
	63.21	77.3	94.3	230.01	370.1	531.02	701.15
	63.3	77.4	94.5	230.03	370.2	531.03	701.24
	63.4	78.11	94.9	230.04	371.02	532.01	701.25
	63.53	78.12	95.2	232	371.11	532.02	701.33
	63.6	78.3	95.9	233	371.21	533.01	703.12
	63.7	79.3	97.4	234	372	533.02	703.13
	63.81	79.4	97.5	235	373.03	533.03	703.22
	63.82	79.5	98	236	373.04	534.01	705
	63.91	81.2		237	373.11	534.02	901.03
	63.92	81.3		238	373.21	535.1	901.22
	64.1	81.4		240.02	416.05	535.2	902.02
	67.1	81.6		241.01	417.02	536.02	
	68.21	82.1		241.02	433.3	537.22	
	69.31	82.4		242	434.02	538.11	
	69.41	82.91		243	436.13	538.12	
	69.44	83.11		244.01	436.23	539	
	69.45	83.12		244.22	437.11	540.01	
	69.5	83.22		245.12	437.12	540.01	
	69.9	83.3		245.22	437.22	540.22	
	71.11	83.4		247.2	437.32	541.2	
	71.12	83.5		248	438.06	542.02	
	71.13	83.6		249.03	438.21	542.11	
	71.2	82.7		249.22	438.31	542.97	
	72	83.8		249.32	440.06	543	
	73.9	83.91		250	441.02	544	
	74.1	85		251	444.04	545.01	

## **IV. DATA COLLECTION METHODOLOGY FOR THE COMPARISON GROUP**

### **4.1 Description and Schedule**

The data collection for the comparison group for this study was conducted from Mathematica's Telephone Center in Princeton, New Jersey, using computer-assisted telephone interviewing (CATI). The interviewer training took place on August 1, 1992. During September 1992, a decision was made to increase the sample size from 1,500 smokers to 1,500 households. The final interviews were completed by October 31, 1992.

An adult member of the household who was 18 years of age or older was interviewed about the characteristics of all smokers 12 years of age and older in the household. One minute per call was spent attempting to establish any contact, whether answered or not, with a sampled household. After an eligible member of a household with at least one smoker was contacted, four minutes were spent on average to complete an interview.

Table IV.1 depicts the distribution of smokers in the sample. Self-reporters are respondents who smoke and provided information for themselves. A smoker proxy is a smoker who provided information for another smoker in the household. A nonsmoker proxy is a nonsmoker who provided information for a smoker in the household. As the table indicates, 51.1 percent were self-reporters, 26.2 percent were smoker proxies, and 22.7 percent were nonsmoker proxies.

Table IV.2 represents the final "closeout" status of all households originally included in the sample. Households that were outside the fire service areas or that provided incomplete or incorrect information that made matching the survey data to the manufacturer's data on cigarette characteristics impossible were deleted or flagged. As Table IV.2 shows, an average of 3.9 households had to be screened to reach a household with at least one smoker. In other words, a little over one quarter of the sample of households had a smoker. The next section provides details on response rates. Appendix C provides details on the final status of cases sampled for each of the eight sites.

TABLE IV.1

DISTRIBUTION OF SMOKERS IN THE SAMPLE BY TYPE OF INTERVIEW

Type	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Self-reported	1,128	51.1	1,128	51.1
Smoker proxy	577	26.2	1,705	77.3
Nonsmoker proxy	501	22.7	2,206	100.0

TABLE IV.2

## FINAL STATUS OF HOUSEHOLDS ORIGINALLY SAMPLED--TOTALS

	Total Households		
	Number	Percent	Calls
<b>Eligible</b>			
Complete	1,503	12.9	2.9
Complete--no address	29	0.2	5.8
Final refusal	65	0.6	6.1
Subtotal	1,597	13.7	3.1
<b>Ineligible Residence</b>			
No smoker > 12	4,618	39.7	3.0
Subtotal	4,618	39.7	3.0
<b>Eligibility Unknown</b>			
Language barrier	24	0.2	3.8
Final refusal	385	3.3	5.4
Maximum dialings	567	4.9	20.0
Effort ended	91	0.8	12.3
Other	1	0.0	18.0
Subtotal	1,068	9.2	13.7
<b>Nonresidence</b>			
Nonworking/new number	2,317	19.9	1.5
Not a residence	2,036	17.5	2.3
Other	3	0.0	2.7
Subtotal	4,356	37.4	1.9
<b>Total</b>	<b>11,639</b>	<b>100.0</b>	<b>3.6</b>



## 4.2 Response Rates Defined

In this section we consider two important measures of survey quality for this study: the overall response rate ( $RR_{\text{overall}}$ ), and the response rate for smokers ( $RR_{\text{smokers}}$ ). A response rate is the ratio of the number of completed interviews with reporting units to the number of eligible reporting units in the sample. (For computing  $RR_{\text{overall}}$  an eligible reporting unit is a telephone household. For computing  $RR_{\text{smokers}}$  and eligible reporting unit is a telephone household with one or more smokers.) A response rate is simple to compute when the eligibility status of every reporting unit in the sample is known. When eligibility is not known, assumptions need to be made about how many of the "unknown" units are eligible.

In this study, the eligibility status of every unit is not known. For instance, some households may have refused to be interviewed before we could establish whether the households contained a smoker. Some telephone numbers were retired after a maximum number of dialings, without ever making a contact that would allow us to determine whether the number served a household. In these cases, we used the results of the survey to estimate two rates: a rate of smoking in known households, and a rate of household numbers among the sample telephone numbers. We used these survey-based rates to estimate the total number of households and the total number of households with one or more smokers.

The overall response rate is the number of households that completed the screening part of the interview ( $C_{\text{screen}}$ ) divided by the estimated number of households in the sample. Both smoking and nonsmoking households completed the screening part of the interview. The denominator for the overall response rate is the number of known households (HH), plus an estimate of the number of households called without successfully determining whether or not the number belonged to a household ( $N_{\text{HHDK}}$ ). This estimate uses a so-called "household rate" (HHRATE), which is an estimate of the proportion of all telephone numbers in the sample that are household telephone numbers. Thus, the overall response rate is:

$$RR_{OVERALL} = \frac{C_{SCREEN}}{HH + (N_{HDK} * HHRATE)}$$

with HHRATE is defined as:

$$HHRATE = \frac{HH}{HH + N_{NR}}$$

where HH is the number of households identified and  $N_{NR}$  is the number of nonresidential numbers identified.

The smoker response rate is the number of completed interviews with smoking households ( $C_s$ ) divided by the estimated number of smoking households ( $HH_s$ ). The denominator of the smoker response rate is the sum of: (1) the number of households known to have smokers ( $HH_{SK}$ ); plus (2) a survey-based estimate ( $HH_{SE}$ ) of the number of other households containing smokers. Thus:

$$RR_{smoker} = \frac{C_s}{HH_s} = \frac{C_s}{HH_{SK} + HH_{SE}}$$

The estimate  $HH_{SE}$  has two components: (1) a portion of the telephone numbers known to be households, but where smoking status was unknown ( $HH_{SDK}$ ), and (2) a portion of the numbers called without determining whether or not they were household numbers ( $N_{HDK}$ ). Thus:

$$HH_{SE} = (HH_{SDK} * SRATE) + (N_{HDK} * HHRATE * SRATE)$$

For SRATE, we used the prevalence of smokers among households, as estimated from the survey:

$$SRATE = \frac{HH_{SK}}{HH_{SK} + HH_{NON}}$$

where  $HH_{NON}$  is the number of households known not to have smokers.  $HHRATE$  is the same rate used in the overall response rate, discussed earlier.

For the entire sample, the overall response rate was 87.17%. The response rate for smokers was 83.30%. Tables IV.3 and IV.4 summarize the response rates by site, by income strata, and for all sites combined, respectively.

Neither the overall response rate nor the response rate for smokers varied greatly between fire service areas. (The lowest rates were in Philadelphia.) The overall response rate ranged from 85.69% to 89.05%, and the response rate for smokers ranged from 80.64% to 87.18%.

TABLE IV.3

RESPONSE RATES BY SITE

	Baltimore	Cleveland	Columbus	Dallas	Denver	Houston	Philadelphia	Portland
RR <sub>smokers</sub>	82.95	85.77	87.18	84.33	83.54	82.86	80.64	81.89
RR <sub>overall</sub>	87.43	89.05	88.78	87.84	87.19	86.76	85.68	85.94
SRATE	30.74	33.18	24.42	23.67	23.37	22.92	27.21	25.86
HHRATE	65.29	55.77	67.79	54.42	51.37	55.49	74.59	60.82
Smoker Households	157	141	167	225	107	317	311	107

TABLE IV.4

## RESPONSE RATES BY INCOME STRATUM

	Low Income Exchanges	Other Exchanges	All Exchanges Combined
RR <sub>smokers</sub>	84.62	83.01	83.30
RR <sub>overall</sub>	87.88	87.02	87.17
SRATE	31.19	24.61	25.69
HHRATE	46.41	64.09	60.32
Smoker Households	306	1,226	1,532

## V. DATA PROCESSING AND ENTRY

Data collected by the CATI system require little editing or coding. The "other, specify" answers for brand codes were printed out. Some answers were matched to existing codes. The remaining cigarette brands given under "other, specify" are listed in Appendix D. Frequency distributions on the "cleaned" survey data do not indicate any unexpected or unreasonable values.

## VI. ADDITION OF CENSUS AND OTHER DATA

### 6.1 Census Tracts

Census tracts were identified for all but 61 interviewed households. The file layout shows a "source of census tract" code in column 149. There were three sources for the census tract.

If the randomly generated sample telephone number belonged to a listed telephone household, then a census tract was added to the sample record on the basis of the address published with the listed telephone number. These cases are identified by a 1 in column 149.

When the telephone number did not belong to a listed telephone household, the address supplied by the respondent in the interview was used to identify a census tract. (The respondent-supplied address is found in columns 68-148). This method used computer matching of the respondent address with a file containing street names, and house number ranges from each ZIP code and census tract. Census tracts identified through this method are denoted by a 2 in column 149.<sup>1</sup>

In cases where the computer match failed, a census tract was imputed exclusively on the basis of the ZIP code provided by the respondent. This imputation added the census tract associated with the ZIP code's center of population. A value of 3 in column 149 indicates that the tract was imputed.<sup>2</sup>

---

<sup>1</sup>In 95 cases, the household address could not be matched with a census tract. Failure to match may have occurred for several reasons: The respondent may have misreported the street name or house number (perhaps thinking that deliberate misreporting would ensure confidentiality), or this information may have been misrecorded by the interviewer. Respondents who refused to give an address were asked to provide an intersection near their house. Some of these "intersections" were found to be nonexistent.

<sup>2</sup>The census tracts imputed on the basis of ZIP code should not be relied on with the same confidence as the census tracts obtained using the other methods. An urban ZIP code area can contain many census tracts (we estimate an average of nine tracts per ZIP code), so values of 3 in column 149 should be interpreted as a warning that the census tract is only an approximation. Values of 1 and 2 indicate that the census tract can be relied on to be correct.

## 6.2 Census Data

Three pieces of tract-level 1990 census data were added to each record in the final version of the survey file. They were: (1) median household income for the tract; (2) the percentage of the population aged 25 and above who had completed at least a high school diploma; and (3) percentage of persons in the tract below the poverty level. These data were obtained for each tract in the eight cities. The quality of the data was verified against paper copies of the same data. No errors or omissions were found.

## 6.3 Other Merged Data

The survey data for a given case in the sample was matched to the cigarette characteristic data provided by cigarette manufacturers. The cigarette characteristic data included information on the following: density, porosity, citrate, and circumference (from which the amount of tobacco could also be calculated). The data from the two sources were matched on a seven-digit code called a "key code." The seven digits of this code are as follows:

1st digit	Manufacturer
2nd and 3rd digits	Brand code
4th digit	Length of cigarette
5th digit	Filter or not
6th digit	Soft or hard pack
7th digit	Mentholated or not

Of the 2,206 smokers in the original sample for the comparison group, 1,969 were matched to data provided by cigarette manufacturers.

The following rules applied for matching and merging the two data sets (the cigarette characteristic data provided by the manufacturers will be referred to as CCD):



- (1) If the survey UPC matched a CCD UPC and the survey key matched a CCD key then, if the brand codes were the same, the UPC match was used. Otherwise the key match was used.
- (2) If the survey UPC matched a CCD UPC but there was not match on keys, the UPC match was used.
- (3) If there was no match on UPC codes but there was a match on keys, the key match was used.
- (4) If there was one missing element (other than brand code) in the survey key and there was only one match with the CCD when this element was excluded, the corresponding CCD key was used.
- (5) If there were no missing elements in the survey key and that key matched only one CCD key when only one element was excluded, the corresponding CCD key was used.
- (6) If a match as in the previous two items resulted in more than one potential CCD key match but the characteristics were identical across all the potential matches then one of the potential matches was arbitrarily selected.

Table VI.1 depicts the status of the cases in the sample when matched to the manufacturer's data.

TABLE VI.1

NUMBERS AND PERCENTAGE OF SURVEY CASES  
MATCHING MANUFACTURER'S DATA

	Numbers	Percentage
Matched by UPC	717	32.5
Matched by key code	1,059	48.0
Inferred match	193	8.7
No match possible due to missing data	200	9.1
No match possible due to no available key code in manufacturer's data set	37	1.7

## VII. USING THE DATA

The data will be used to make descriptive statements about smokers in the targeted fire service areas and to compare smokers associated with household fires with smokers not associated with fires. Data on smokers associated with fires were collected in a separate study.<sup>1</sup> Logistic regression models will be used to determine the effect of various characteristics of smokers and of cigarettes on the probability that a household fire occurs.

In this section we suggest guidelines for using the data collected in the survey. Subsection 7.1 describes the limitations of the data. Subsections 7.2 and 7.3 deal with the specific issues of weighting, imputation and computing sampling error.

### 7.1 Limitations in Using the Data

Data from this survey are subject to the usual limitations of survey data. The data are affected by several sources of potential error:

- Sampling error because the data were collected from a sample of smokers, rather than the entire population
- Error arising from non-response (both case and item level), and possible frame undercoverage or overcoverage
- Response error due to questions being misinterpreted or information incorrectly recalled by respondents
- Interviewer or processing error

In addition to the general issues of data limitations, there are conditions present in this survey that affect the usefulness of the data for comparison with the data collected in households with fires. While

---

<sup>1</sup>In the present survey nine smokers were identified whose households had experienced fires. We suggest excluding these households from analyses that combine data from the two studies. The number of fires observed in the present survey is very small and these smokers should have had a chance of being included in the household fire study.

these conditions do not invalidate the comparisons, they should be considered as possible sources of "noise" in conducting the analyses. First, although both samples were drawn from the same fire service areas the actual coverage is somewhat different. The present survey collected data only from smokers in telephone households, while the study of households with fires collected data from smokers thought to cause a fire and is without regard to presence of a telephone. Second, as is usually the case, there may be method effects. The present survey was conducted by telephone on all smokers in a household, with one respondent reporting for all smokers in the household therefore, proxy data is collected. The household fire study collected data only for the smoker in a household who was suspected of causing a fire, used in-person interviewing, and allowed for proxy responses only in cases when the desired respondent was unavailable due to injury. However, the most recent literature indicates little difference in data quality between telephone and inperson methodologies.<sup>2</sup>

Despite these potential concerns, there is little reason to expect substantial biases from the use of different data sources. The vast majority (over 90 percent) of households in these areas have telephones. The degree of geographic undercoverage was small, and the identification of the census tracts for over 96 percent of the sample insures that few persons outside the service areas will be included in the final data set. Given these considerations, it is our opinion that the resulting bias will be small, but we cannot directly measure the extent of the bias from the survey data.

**Sampling Error.** The sample for this survey is not a simple random sample, and therefore proper analysis of the data requires that the effects of departures from simple random sampling (called design

---

<sup>2</sup>De Leeuw, Edith D., and Johannes van der Zouwen. "Data Quality in Telephone and Face to Face Surveys: A Comparative Meta-Analysis." In *Telephone Survey Methodology*, edited by Paul P. Biemer, Robert M. Groves, Lars E. Lyberg, James T. Massey, William L. Nicholls II, and Joseph Wakesberg. New York: John Wiley & Sons, Inc., 1988, pp. 283-299.

effects)<sup>3</sup> be assessed and taken into account in conducting analyses, and interpreting and presenting results. Section 7.3 below suggests methods of computing sampling error.

**Potential Bias Due to Non-Response and Coverage Problems.** Error from non-response and from undercoverage arise because potential respondents could not be interviewed or were excluded from the sample frame, respectively and the omitted individuals may differ substantially from those that were interviewed. Overcoverage means that some persons living outside the service areas may have been interviewed. Overcoverage can cause bias if those who are erroneously included differ from the study population.

The major sources of undercoverage (discussed in more detail in Section 3.2) relates to selecting telephone exchanges for the RDD sample. For reasons of efficiency, we excluded exchanges in which only a small proportion of the exchanges' listed telephone numbers were in the service area. However, the excluded exchanges never contained more than 5 percent of the total telephone households in the service area.

The high response rate achieved (83% for smokers, see Section 4.2) alleviates much of the concern regarding non-response. Steps taken to identify interviewed households living outside the service area reduce the impact of overcoverage. The strategy of defining strata by household income and allocating the sample proportional to all households should offset the problems related to undercoverage.

**Response, Interviewer and Processing Efforts.** Response error, and interviewer and processing errors, occur when the respondent (intentionally or unintentionally) gives incorrect reports, or when the responses are incorrectly coded or changed in processing. The use of computer-assisted interviewing

---

<sup>3</sup>A design effect (Deff) is the ratio for the sample variance given the actual sample design to the variance that would be obtained with a (hypothetical) simple random sample (SRS) of the same size. Thus:

$$Deff = \frac{Var(DSIGN)}{Var(SRS)}$$

SRS estimates of standard errors are multiplied by the square root of the design effect (deft) to obtain more accurate estimates of standard errors for constructing confidence intervals or performing significance tests.

greatly reduces interviewer and processing errors. Interviewer errors on this study should also be small as a result of the extensive training conducted. Frequency distributions of the data do not indicate any unexpected or unusual values.

## **7.2 Need for Weighting and Imputation**

Sample weights are used when the sample is distributed differently, on important characteristics, than the study population. The differences in distribution may result from the study design (e.g., oversampling) or from differential response rates or frame coverage. Imputation refers to a set of procedures for adding values for items missing from cases that are otherwise complete.

**Weights.** The choices of whether to use sample weights, and if so, how to construct the weights, depends on the use of the data. In the present survey, the data may be used separately to make descriptive statements about smokers in the fire service areas. Their primary use will be, combined with data from the household fire survey, to analyze the effects of smoker and cigarette characteristics on the likelihood of a household fire occurring.

With regard to making descriptive statements, the sample was designed to provide estimates that are self-weighting with respect to all smokers in the areas included in the sample frame. Sample strata were defined by fire service area and median household income. By targeting the sample so that completed interviews with eligible households are distributed across strata in approximately the same proportion as are estimates of all households, the need to use weights for descriptive analysis should be eliminated. This approach to design has provided the advantage of explicitly controlling for the distribution of households across low-income and other areas. If we were to compute sample weights, we would calculate the stratum-specific probabilities of selection and response rates, and weight by the inverse of these. We would then check the weighted sample distribution against our best estimate of the population distribution and adjust to that distribution.<sup>4</sup> However, since the sample distribution already fits that of

---

<sup>4</sup>We would prefer to use external estimates of the distribution of smokers, but since we do not have such estimates, the distribution of households allows us to estimate the distribution of smokers from

the estimated population, we have accomplished by sample control, what would have been done by weighting. Table VII.1 indicates by stratum the estimated distribution of households the population, and cooperative<sup>3</sup> households identified in the sample. Based on the distribution obtained, it is not necessary to use sample weights in conducting the data analysis.

While we conclude that sample weights are not needed, there are other approaches to descriptive analysis that would lead to different decisions about weighting. For instance, to make estimates about smokers in telephone households one would weight to the estimated distribution of all telephone households in the frame or in the service areas.

For regression analysis using only the present survey, weights should not be required, even if one took a different view toward weighting for descriptive analysis. Any household weights would be constant within strata. In a multiple regression model, variables can be added to control for the stratum-specific effects that would be addressed by the sample weights that would be computed.

For analyses combining data from the present survey with that from the household fire survey, weights should not be needed if the objective is to estimate the coefficients of the independent variables in a logistic regression model. However, to estimate the likelihood of a household fire, whether unconditionally, or conditioned on certain values of the dependent variables, weights would

---

the sample. The prevalence of smokers in cooperative households provides our best estimate of the prevalence of smokers in the population.

<sup>3</sup>By cooperative households we mean those that provided information on the number of smokers in the household.

TABLE VII.1

## DISTRIBUTION OF HOUSEHOLDS IN SAMPLE AND POPULATION BY STRATUM

Stratum	Percent of Sample Households in Stratum	Estimated Percent of Population Households in Stratum
Baltimore Low-Income	1.10	1.16
Baltimore Remainder	7.59	7.45
Cleveland Low-Income	2.64	2.82
Cleveland Remainder	4.46	4.99
Columbus Low-Income	1.28	1.30
Columbus Remainder	9.88	9.70
Dallas Very Low-Income	0.33	0.42
Dallas Low-Income	1.66	1.83
Dallas Remainder	13.72	13.46
Denver Low-Income	1.04	1.15
Denver Remainder	6.64	6.59
Houston Low-Income	3.69	3.74
Houston Remainder	19.54	19.15
Philadelphia Low-Income	4.10	3.94
Philadelphia Remainder	15.40	15.18
Portland Low-Income	0.58	0.61
Portland Remainder	6.34	6.51
Total	100.00	100.00



be required, because households experiencing fires would be substantially overrepresented in the combined data set.<sup>6</sup> In such a case, weights should be constructed to reflect the distribution of smoking households by fire/non-fire status.

**Imputation.** Imputation is used to compensate for missing items within otherwise complete interviews. We do not recommend replacing any missing items (especially keycode items) with imputed values on a *permanent* basis. If values were imputed, the data set would have less precision than a data set of the same size with no item non-response. Thus, standard deviations calculated on the imputed data would be underestimated and other descriptive statistics may be distorted. Further, the use of artificial values would make matching the data with the manufacturer's data imprecise.

In conducting regression analysis, however, to avoid dropping a large number of cases, a procedure may be employed that imputes values where cases are missing data for independent variables other than those included in the "keycode." In this procedure, a constant value (usually zero or the sample mean) is imputed for the missing variable(s). A binary variable is then created for each variable where values are imputed. For each case, the binary indicator is set to 1 if the variable was originally missing and set to zero otherwise. The binary indicator for a variable is then included in any regression equation that contains that variable.

### 7.3 Computing Sampling Errors

The effects of departures from simple random sampling are usually grouped as the design effects of clustering ( $Deff_c$ ) and weighting ( $Deff_w$ ). Although we recommend that sample weights not be used in the analysis, we realize that some approaches to analyzing the data could call for weights. Thus, we will briefly explain why weights affect sampling error, and how one would estimate  $Deff_w$ , should that be required by some future part of the analysis.

---

<sup>6</sup>In a regression model, not using weights would bias the intercept term but not the coefficients of the independent variables.

The design effects of weighting results from the use of weights to compensate for differential sampling rates and non-response. A weighted estimate (e.g., a mean) is not a simple statistic but a complex one involving two variables--the variable of interest and the weighting variable. Thus, the estimated variance of a weighted statistic must account for two sources of variation.<sup>7</sup>

The design effect of clustering  $Deff_c$  reflects the fact that in a clustered sample, the units being observed are not selected independently, but as part of larger units known as clusters.<sup>8</sup> The variance of an estimate from a clustered sample has two components--between clusters and within clusters.<sup>9</sup> In the present survey, the household is the cluster and individual smokers the unit of observation. Clustering, like weighting can affect the sampling error of any statistic.

---

<sup>7</sup>The effect of weighting on sampling error can be estimated using the methods described below for estimating clustering effects. A useful approximation for the design effect of weighting ( $Deff_w$ ) is 1 plus the relvariance ( $rv$ ) of the weights:

$$Deff_w = \frac{Var(Weighted)}{Var(SRS)} \doteq 1 + rv$$

$$rv = \frac{n s_w^2}{\sum W_i}$$

where:

$n$  is the (unweighted) number of cases

$$s_w^2 = \left[ \frac{1}{n-1} \right] \left[ \sum_{i=1}^n w_i^2 - \frac{\left[ \sum_{i=1}^n w_i \right]^2}{n} \right]$$

$w_i$  is the weight for the  $i$ th case.

<sup>8</sup>Clusters are called primary sampling units (PSUs) when there is more than one stage of sampling. In the present survey, we sampled all eligible persons in a household, so using the term PSU to refer to households could be confusing.

<sup>9</sup>In estimating sampling error, the within cluster component of variance would be zero, since all smokers in a household were sampled.

There are several methods for estimating the standard errors for statistics from a complex sample. These methods we will discuss fall into two groups: Taylor series approximations and replication or resampling methods.<sup>10</sup> For the present survey, the most important statistics are regression coefficients, for which replication methods (e.g., jackknife, balanced repeated replications) are usually preferred.

To estimate the variance of means or proportions for smokers, one could use commercially available packages, such as WESVAR or SUDAAN, or one could use SAS or SPSS to estimate the components of a variance estimate for a ratio mean based on the Taylor series. (MPR has written SAS routines to perform these computations.) In the Taylor series approximation, we define:<sup>11</sup>

$a_h$  = the number of households in stratum  $h$

$x_{ah}$  = the number of smokers in the  $a$ th household in the  $h$ th stratum

$y_{ah}$  = the value of the variable  $y$  for the  $a$ th household in the  $h$ th stratum

$x$  = the total number of cases (smokers) =  $\sum_{h=1}^H \sum_a x_{ah}$

$y$  =  $\sum_{h=1}^H \sum_a y_{ah}$  --(the sum of variable  $y$  across all households in all strata)

$r$  =  $y/x$  (the ratio mean)

$H$  = the number of strata

---

<sup>10</sup>Other groups of methods include generalized variance functions and random group methods. All these methods are explained in Wolter (1985). Kalton (1983) gives easy to understand examples of some of these methods.

<sup>11</sup>The formulae below were taken from Kalton, (1983), pp. 44-45. Equivalent formulae are found in Kish (1965), p.192. The formula for  $V(r)$  is equivalent to that found in Wolter (1985), p.236.

$$v(y) = \sum_h a_h s^2(y)_h$$

$$v(x) = \sum_h a_h s^2(x)_h$$

$$c(x,y) = \sum_h a_h s(x,y)_h$$

where:

$$s^2(y)_h = \sum_a [y_{ah} - (\sum_a y_{ah}/a_h)]^2 / (a_h - 1) \text{ --(the sum of the within stratum variances of y)}$$

$$s^2(x)_h = \sum_a [x_{ah} - (\sum_a x_{ah}/a_h)]^2 / (a_h - 1) \text{ --(the sum of the within stratum variances of x)}$$

$$s(x,y)_h = \sum_a [x_{ah} - (\sum_a x_{ah}/a_h)] [y_{ah} - (\sum_a y_{ah}/a_h)] / (a_h - 1) \text{ --(the sum of the within stratum covariances of x and y)}$$

The variance  $v(r)$  of the ratio mean  $r$  is then approximately:

$$v(r) \doteq [v(y) + r^2 v(x) - 2r (c(x,y))]/x^2$$

The standard error of the ratio mean would be the square root of  $v(r)$ . Standard errors can be estimated for all variables of interest, or estimates of an average design effect can be calculated. Simple random sample estimates of standard errors are then multiplied by the square root of the average design effects.

A less precise, but useful approximation, since no subsampling was done within households, would be to compute standard errors as if the household were the unit of observation. Thus for a statistic  $y$ :

$$v(y) \doteq \sum_{h=1}^H w_h^2 \frac{s^2(y)_h}{a_h}$$

where:

$w_h$  is the proportion of the population in stratum  $H$ .

For regression coefficients, we recommend a jackknife or balanced repeated replications (BRR) approach. These are available for logistic regression in WESLOG or CPLEX. SUDAAN will compute standard errors for logistic regression coefficients, but uses the Taylor series approximations. CPLEX is available free of charge from the Bureau of the Census. SUDAAN is available from Research Triangle Institute and WESLOG from Westat, Inc. MPR has SUDAAN and is obtaining CPLEX.

In a jackknife estimation of the variance of regression coefficients, the sample is divided into  $k$  random groups, each of size  $m$ . In the present case we would divide the households, rather than the smokers, into groups. One could use each household as a group or could form larger groups. The jackknife estimate of standard errors requires  $k + 1$  estimates of the coefficients, one with all cases in the model, plus  $k$  estimates, each with one "group" omitted. Then:

$\hat{B}$  = the regression coefficient with all cases in the model

$\hat{B}_a$  = the regression coefficient with the  $a$ th random subgroup omitted

$$B_a^* = \frac{1}{k} \sum_{a=1}^k B_a$$

The variance of  $B_a^*$  is then:

$$v(B_a^*) = \frac{1}{k(k-1)} \sum_{a=1}^k (\hat{B}_a - \hat{B})^2$$

One then may use  $v(B_a^*)$  directly or compute average design effects, where:

$$Deff_c = \frac{v(B_a^*)}{v(\hat{B})}$$

## VIII. FREQUENCY AND CROSS-TABULATIONS OF THE SURVEY DATA

Frequencies, means for continuous variables, and crosstabulations of the survey data for the comparison group are provided in Appendix E.

Frequency distributions are provided for the following in the Appendix, listed in order of appearance:

- Number of smokers in a household
- Type of respondent
- Gender of smoker
- Length of cigarette
- Filtered, nonfiltered
- Package type
- Mentholated or not
- Number of cigarettes smoked (amount)
- Age of smoker
- Race of smoker
- Hispanic or not
- Education of smoker
- Persons living in household
- Household income
- Ownership of home
- Presence of cigarette fire
- Percentage below poverty (census tract)
- Median income (census tract)
- Percentage with high school education (census tract)

Means, medians, and the minimum and maximum values were computed for:

- Number of smokers in a household
- Age of smokers
- Number of persons in a household
- Percentage below poverty in tract
- Median income for tract
- Percentage with high school education in tract

Cross-tabulations are provided for the following, listed in order of appearance:

Age \* Gender  
Age \* Race  
Age \* Education  
Age \* Percentage with high school education in tract  
Age \* Income  
Age \* Median income in tract  
Age \* Percentage below poverty in tract  
Age \* Filtered cigarette  
Age \* Number of cigarettes (amount)  
Age \* Density  
Age \* Amount of tobacco  
Age \* Porosity  
Age \* Citrate

Race \* Gender  
Race \* Education  
Race \* Percentage with high school education in tract  
Race \* Income  
Race \* Median income in tract  
Race \* Percentage below poverty in tract  
Race \* Filtered cigarette  
Race \* Mentholated  
Race \* Number of cigarettes (amount)  
Race \* Density  
Race \* Porosity  
Race \* Citrate

Gender \* Education  
Gender \* Percentage with high school education in tract  
Gender \* Income  
Gender \* Median income in tract  
Gender \* Percentage below poverty  
Gender \* Filtered cigarette  
Gender \* Number of cigarettes (amount)  
Gender \* Density  
Gender \* Porosity  
Gender \* Citrate

Income \* Median income in tract  
Income \* Percentage below poverty  
Income \* Filtered cigarette  
Income \* Number of cigarettes (amount)  
Income \* Porosity  
Income \* Citrate

Education \* Percentage with high school education in tract  
Education \* Filtered cigarette  
Education \* Number of cigarettes (amount)  
Education \* Porosity  
Education \* Citrate

Density \* Filtered cigarette  
Density \* Circumference  
Density \* Porosity  
Density \* Citrate

Porosity \* Filtered cigarette  
Porosity \* Circumference  
Porosity \* Citrate

Citrate \* Filtered cigarette  
Citrate \* Circumference

Age \* Amount of tobacco  
Race \* Amount of tobacco  
Gender \* Amount of tobacco  
Income \* Amount of tobacco  
Education \* Amount of tobacco  
Density \* Amount of tobacco  
Porosity \* Amount of tobacco  
Citrate \* Amount of tobacco



Contract No.: CPSC-C-92-1001, Task Order 005  
MPR Reference No.: 8071

**SELF-PROXY COMPARISONS FOR THE  
CIGARETTE FIRE SAFETY SURVEY**

**FINAL REPORT**

OMB #3041-0110

February, 1993

Authors:

Donna Eisenhower  
John Hall  
Randy Brown

Submitted to:

U.S. Consumer Product Safety Commission  
5401 Westbard Avenue  
Bethesda, Maryland 20816

Submitted by:

Mathematica Policy Research, Inc.  
P.O. Box 2393  
Princeton, N.J. 08543-2393  
(609) 799-3535

Project Officer:

William Zamula

Project Director:

Donna Eisenhower

## CONTENTS

Chapter		Page
I	BACKGROUND FOR THE STUDY .....	B-1
II	PRELIMINARY COMPARISONS OF SELF AND PROXY-REPORTED DATA FROM THE ORIGINAL SURVEY .....	B-5
III	DESIGN OF THE REINTERVIEW STUDY .....	B-13
IV	FINDINGS BASED ON THE REINTERVIEW DATA .....	B-15
V	CONCLUSIONS .....	B-31

## TABLES

Table		Page
II.1	ORIGINAL SURVEY SMOKER BEHAVIOR AND CIGARETTE CHARACTERISTICS . . . . .	B-6
II.2	SMOKER BEHAVIOR AND CIGARETTE CHARACTERISTICS BY SEX, HISPANIC, RACE AND HOME OWNERSHIP . . . . .	B-8
II.3	SMOKER BEHAVIOR AND CIGARETTE CHARACTERISTICS CONTROLLING FOR SEX, RACE AND HOME OWNERSHIP . . . . .	B-10
IV.1	PERCENTAGE OF CASES WHERE REINTERVIEW RESPONSES DO NOT MATCH ORIGINAL RESPONSE, BY TYPE OF RESPONDENT TO ORIGINAL INTERVIEW . . . . .	B-16
IV.2	NUMBER OF DAYS ELAPSING BETWEEN ORIGINAL INTERVIEW AND REINTERVIEW BY TYPE . . . . .	B-18
IV.3	FREQUENCY OF BRAND CHANGE BY TYPE OF REPORT AND BRAND CODE MATCH STATUS . . . . .	B-20
IV.4	A COMPARISON OF THE DISTRIBUTION ON LENGTH OF CIGARETTE FOR THE ORIGINAL AND REINTERVIEW SAMPLE . . . . .	B-22
IV.5	A COMPARISON OF THE DISTRIBUTIONS OF PACK TYPE FOR THE ORIGINAL AND REINTERVIEW SAMPLE . . . . .	B-23
IV.6	NON-MATCH RESPONSES TO INCOME QUESTION BY VERSION, OVERALL AND BY TYPE OF ORIGINAL RESPONDENT . . . . .	B-25
IV.7	MISSING DATA BY TYPE OF ORIGINAL RESPONDENT . . . . .	B-27
IV.8	NON-MATCH RESPONSES TO INCOME QUESTION BY VERSION, OVERALL AND BY TYPE OF ORIGINAL RESPONDENT . . . . .	B-29

## I. BACKGROUND FOR THE STUDY

On August 10, 1990, Congress passed The Fire Safe Cigarette Act of 1990. The act authorized the U.S. Consumer Product Safety Commission (CPSC) to conduct research and assess the feasibility of developing a performance standard to reduce cigarette ignition propensity. Data have now been collected by two organizations which will help the CPSC determine the relationship between various characteristics of cigarettes and smokers and the risk of fire.

The National Fire Protection Association (NFPA), under contract with the Consumer Product Safety Commission, has undertaken a fire-incident study, based on data they have collected on cigarette-related fires at eight sites. The data collection began in November 1991 and was completed in December 1992. Personnel of participating fire departments were trained to collect the information in person at the scene of a fire. Mathematica Policy Research, Inc., under subcontract with Market Facts, Inc., has completed the collection of data for a comparison group, to be used in determining the effect of characteristics of smokers and cigarettes on the probabilities of household fires.

Mathematica was also responsible for the design, collection, and analysis of this methodological study to evaluate the data quality of self and proxy reports used in the original Comparison Survey. This study was done by comparing results between the original respondent whether a self-report, smoker proxy report, or non smoker proxy report to actual self-reports at a reinterview. The methodological study was conducted in response to concerns expressed by the Technical Advisory Group created by the Fire Safe Cigarette Act of 1990.

Data for the comparison group was collected for all smokers in a household. The information was reported by one household member 18 years of age or older. For the total of 2,206 smokers, 51.1% were self-reports, 26.2% were smoker proxy reports and 22.7% were nonsmoker-proxy reports. Self-reports for all smokers in a household or selecting one smoker per household in the same numbers was not feasible in the survey. This methodological study assesses the quality of the proxy-reported data.

The proxy information provides a means of obtaining data on more smokers in the fire service area. An important issue is whether the proxy report is as accurate and reliable as the data that would have been obtained from the actual smoker. Self-reported data are usually assumed to be more accurate and reliable. However, the survey literature suggests that the distribution of responses from proxies often differs from that of self-respondents without allowing us to conclude which is better. This is because there is rarely an external means available or used to validate the self and proxy reports, or the study design is limited in some other manner.

Moore (1988) after completing a review of the literature on self-proxy reporting spanning three decades concludes that this "research has not produced conclusive evidence of consistent response bias or response error variance differences due to the self/proxy status." He attributes this finding to the methodological shortcomings of much of this literature but cautions that "lack of convincing evidence of quality differences is not synonymous with convincing evidence of no quality differences." The literature is further complicated by findings such as these reported by Mathiowetz and Groves (1985) in reviewing the health survey literature, they found that "although early studies indicate less agreement between the interview report and medical record data for proxy reports than for self reports, more recent studies indicate no difference in response error by type of respondent, or suggest that in some cases proxy reports may be more accurate."

Whether the self report is of higher quality than a proxy report will depend upon the individual, their circumstances in relation to the subject matter, and the subject matter itself. Proxy reporting for the mentally impaired or for children has been preferred to no data at all. Proxy reporting in cases where a self-report may be subject to a high level of social desirability or sensitivity might be preferred. However, the best report is one that can be recalled and reported most accurately. The acceptability of who will report must be evaluated in light of this criterion.

**This empirical study evaluates the reliability and degree of missing information for self and proxy reports of cigarette-related information. The study is based on comparisons of original responses given by proxies for a smoker to subsequently obtained responses from the actual smoker. Original self-reports are compared to self-reports in reinterview of the same person as a measure of reliability. This difference in test-retest reliability can then be factored out of proxy-self report comparisons to draw some conclusions about the validity of proxy responses. The real issue is whether self-reports provide any higher quality information than proxy reports when problems of reliability that exist even for the self-reported data are factored out.**

**The remainder of this report is organized as follows:**

- preliminary comparisons of the self-reported and proxy-reported data from the original survey**
- the design of the reinterview study**
- the report of the findings based on the reinterview data**
- final conclusions**

## **II. PRELIMINARY COMPARISONS OF SELF AND PROXY-REPORTED DATA FROM THE ORIGINAL SURVEY**

In the original survey, one respondent in each household answered questions not only about household level data, but about the personal characteristics and smoking behavior of all smokers identified, plus the characteristics of cigarettes smoked by all smokers. Although the respondents were self-selected (interviews were conducted with any adult member of the household 18 years of age or older who either answered the telephone or was the first eligible adult to come to the telephone), it is instructive to see if there are differences in responses by respondent characteristics. For individual level data, respondents are characterized as:

- self-reporters (smokers reporting their own data)
- smoker-proxies (data provided by smokers about other smokers in the household)
- non-smoker proxies (non-smokers providing data on smokers)

This section addresses the differences in the distributions of cigarette-related information as reported by self-reporters, smoker proxies, and nonsmoker proxies in the original survey.

The analysis beginning with Table II.1 consists of a tabular presentation of distributions for various variables. Chi-square ( $X^2$ ) statistics are used to indicate the strength of any differences seen between groups. While the report refers to levels of statistical significance, the size of the percentage difference must be carefully evaluated. Finally, even if the means or distributions are the same for self and proxy-reported data, the proxies may still be reporting differently for individual cases than the smokers themselves would have (with errors balancing out). The reinterview data provide more control for these factors; the results are presented in Section IV.

TABLE II.1

**ORIGINAL SURVEY**  
**SMOKER BEHAVIOR AND CIGARETTE CHARACTERISTICS**  
 (Percent Distribution by Type of Respondent)

	Smoker Self Report	Smoker Proxy	Non-Smoker Proxy	Total Sample
1. Smokes 20 or More Cigarettes a Day	49.8	54.6	44.8	49.9
Sample Size	1,117	535	462	2,114
$\chi^2 = 9.490$ $DF = 2$ $p = 0.009^*$				
2. Soft Pack Cigarettes	71.5	74.7	69.7	71.9
Sample Size	1,092	529	446	2,067
$\chi^2 = 3.125$ $DF = 2$ $p = 0.210^*$				
3. Smokes Menthol	39.3	40.8	38.0	39.42
Sample Size	1,112	552	449	2,093
$\chi^2 = 0.787$ $DF = 2$ $p = 0.675^*$				
4. Smokes Filtered	99.9	96.8	93.2	95.0
Sample Size	1,118	557	456	2,131
$\chi^2 = 6.83$ $DF = 2$ $p = 0.033^*$				
5. Length				
Regular/King	60.7	65.0	77.8	65.3
Long	35.5	31.4	21.7	30.9
Extra Long	4.9	3.7	1.5	3.8
Sample Size	1,114	545	456	2,115
$\chi^2 = 39.22$ $DF = 4$ $p = 0.000^*$				

Table includes only cases where a valid response (other than don't know) was provided

\*p\* is the probability that the  $\chi^2$  statistic would be this large if there were no differences between the groups of respondents. Values of p less than 0.05 indicate statistically significant differences at the percent level.



Among the measures of smoking behavior and cigarette characteristics, differences were found in the amount smoked, whether the smoker uses filtered cigarettes, and the length of the cigarette as noted in Table II.1). Smokers for whom a smoker proxy provided the data are more likely than self-reporters to consume a pack or more a day, but those for whom a non-smoker proxy provided the information are less likely to smoke this much. The explanation for this most likely relates to characteristics of those falling into each group.

Regarding whether filter cigarettes are smoked, and the length of the cigarette, the pattern is more expected. In both cases, the two types of proxy responses (smoker and non-smoker) differ from self-reports in the same direction, with the difference being larger for non-smoker proxy responses. The differences for length of cigarette are quite large (77.8 percent of non-smoker proxies reporting regular or king, compared to 60.7 percent of self reports,) suggesting that this detail is too subtle for many non-smokers to report on accurately. The differences for type of pack and whether the cigarette is menthol were small and not statistically significant.

An analysis was then performed to examine whether differences in reports of smoking behavior and cigarette characteristics could be explained by differences in the types of households or smokers that were reported on.

The first step was to examine differences in smoking behavior and cigarette characteristics by sex and race. The results are presented in Table II.2. Substantial differences in length of cigarettes are found by sex and whether Hispanic. Filtered cigarettes were reported differentially by sex, and to a smaller extent, by race and homeowner status. The number of cigarettes reported smoked differed by sex of smoker, race and whether Hispanic, with differences of 8 to 28 percentage points observed.

The next step was to see if the differences in reports by respondent type remained when personal and household characteristics were controlled. Because race and whether Hispanic overlap, the two categories were combined to include all Hispanics, and three groups of non-Hispanics: White, Black and other. The results of the analysis are shown in Table II.3. When gender is controlled for, the

TABLE II.2

SMOKER BEHAVIOR AND CIGARETTE CHARACTERISTICS  
BY SEX, HISPANIC, RACE AND HOME OWNERSHIP

	Sex		Hispanic		Race			Owns Home	
	Male	Female	Yes	No	White	Black	Other	Yes	No
1. Smokes 20+	54.1	45.5 <sup>a</sup>	24.1	52.4 <sup>a</sup>	57.4	38.0	35.1 <sup>b</sup>	51.7	48.0
2. Filtered Cigarette	93.2	97.1 <sup>c</sup>	96.2	95.0	95.0	96.6	93.0 <sup>c</sup>	94.1	96.2 <sup>a</sup>
3. Length									
Regular or King	73.6	56.1	79.6	63.7	64.1	65.9	69.5	65.2	65.6
Long	24.2	38.4	19.4	32.2	32.2	29.7	27.7	30.5	31.0
Extra Long	2.2	5.6 <sup>c</sup>	1.0	4.1 <sup>c</sup>	3.7	4.4	2.8	4.3	3.4

<sup>a</sup>  $\chi^2$  statistic,  $p \leq 0.10$

<sup>b</sup>  $\chi^2$  statistic,  $p \leq 0.05$

<sup>c</sup>  $\chi^2$  statistic,  $p \leq 0.01$

TABLE II.3

SMOKER BEHAVIOR AND CIGARETTE CHARACTERISTICS CONTROLLING FOR SEX, RACE AND HOME OWNERSHIP

	Self	Smoker Proxy	Non-Smoker Proxy	Total Sample	$\chi^2$	p =
<u>Controlling for Sex</u>						
<u>Smokes 20+ Per Day</u>						
Male	53.4	61.3	47.1	54.1	11.98	0.002
Female	46.7	45.3	41.7	45.5	1.44	0.488
<u>Filtered</u>						
Male	92.8	95.4	91.6	93.2	3.67	0.159
Female	96.7	98.7	96.1	97.1	3.05	0.218
<u>Length</u>						
Male						
Regular/King	69.3	71.7	83.6	73.6	21.7	0.000
Long	27.5	25.8	16.0	24.2		
Extra Long	3.1	2.5	0.4	2.2		
Female						
Regular/King	53.3	55.5	66.1	56.1	9.84	0.043
Long	40.4	39.2	30.6	38.4		
Extra Long	6.3	5.3	3.3	5.6		

TABLE II.3 (continued)

	Self	Smoker Proxy	Non-Smoker Proxy	Total Sample	$\chi^2$	p =
<u>Controlling for Race</u>						
<u>Smokes 20+ Per Day</u>						
Hispanic	18.5	35.3	23.7	24.6	4.78	0.092
White/Non-Hispanic	62.1	59.8	53.0	59.7	6.17	0.046
Black/Non-Hispanic	31.9	49.7	37.4	37.6	13.61	0.001
Other	41.5	63.6	58.8	49.2	2.53	0.282
<u>Filtered</u>						
Hispanic	96.3	98.1	94.4	96.2	0.97	0.617
White/Non-Hispanic	94.9	96.1	92.5	94.8	3.76	0.153
Black/Non-Hispanic	94.8	98.1	95.1	95.7	2.94	0.230
Other	97.6	91.7	82.3	92.9	4.22	0.121
<u>Length</u>						
Hispanic						
Regular/King	70.9	90.2	81.0	79.3	8.84	0.0689
Long	27.8	7.8	19.0	19.7		
Extra Long	1.3	2.0	0.0	1.1		
White/Non-Hispanic						
Regular/King	59.0	61.9	76.2	63.0	24.79	0.000
Long	36.0	39.5	22.9	33.2		
Extra Long	5.0	3.5	0.9	3.9		
Black/Non-Hispanic						
Regular/King	62.7	63.1	73.8	65.4	6.01	0.199
Long	31.6	32.2	23.8	36.1		
Extra Long	5.5	4.7	2.3	4.6		

TABLE II.3 (continued)

	Self	Smoker Proxy	Non-Smoker Proxy	Total Sample	$\chi^2$	p =
Other						
Regular/King	48.8	63.6	68.8	55.9	5.17	0.273
Long	46.3	36.4	18.8	38.2		
Extra Long	4.9	0.0	12.5	5.9		
<b>Controlling for Homeowner</b>						
<b>Smokes 20+ Per Day</b>						
Owns	51.8	56.3	47.0	51.7	4.64	0.098
Other	48.2	52.2	41.1	48.0		
<b>Filtered</b>						
Owns	93.9	96.4	91.9	94.1	4.84	0.089
Other	95.9	97.0	95.5	96.2		
<b>Length</b>						
Owns	60.0	66.3	74.4	65.2	18.70	0.001
Regular/King	34.2	29.7	24.0	30.5		
Long	5.8	4.0	1.6	4.3		
Extra Long						
Other						
Regular/King	61.5	64.5	80.0	65.6	21.00	0.000
Long	34.6	32.1	18.3	31.0		
Extra Long	3.9	3.4	1.7	3.4		

difference in number of cigarettes smoked per day remains for male smokers but not for women. Controlling for gender substantially reduces the difference in the proportion reported to be smoking filter cigarettes; however, substantial differences in length of cigarette remain for both men and women.

When controlling for race, the difference in reports for amount smoked is greater than average among Blacks, and Hispanics and lower among Whites. The pattern of the smoker proxies being more likely than self or non-smoker proxies to report consumption of 20 or more cigarettes a day holds for all the racial groups.

As in the case of controlling for sex, when race is controlled for, differences in reports of smoking filtered cigarettes are greatly reduced. Differences in reports of cigarette length are reduced for Blacks; for Hispanics, the overall pattern changes from non-smoker proxies being most likely to report regular or king size length, to non-smoking proxies being most likely.

Controlling for home ownership reduces the differences on smoking filtered cigarettes, but has little effect on the other two measures.

These comparisons (the usual type of assessment of the validity of proxy responses) suggest that there are sizeable differences between the data reported for smokers who responded to the survey themselves and the data reported for smokers by proxy respondents. Whether these differences are due to reporting error by proxies or to differences between the individuals who responded themselves and those whose information was supplied by a proxy cannot be ascertained from these comparisons. However, differences between the two groups of smokers on basic demographic factors do not appear to account for the differences in the responses. The next section presents a direct assessment of the correspondence between proxy and self reports for the *same individuals*.

### III. DESIGN OF THE REINTERVIEW STUDY

The reinterview sample comprised 600 cases selected from households with three or fewer smokers. First, 200 households were selected where the initial respondent was a non-smoker. Households were selected with probability proportional to the number of smokers, and one smoker was randomly picked for reinterview within each household. Thus, each smoker in the original sample of cases for which a non-smoker proxy provided the data has an equal overall probability of selection for the reinterview sample.

Next, a sample of 200 households was selected from the group where a smoker was the original respondent. Selection was proportional to the total number of smokers minus one. During interviewing a smoker was randomly selected who was not the original respondent.

Finally, a sample of 200 additional households was selected where the smoker was the initial respondent for the household. For this sample, the initial respondent was reinterviewed.

This approach produced 294 completed reinterviews, with 97 that were originally nonsmoker proxy interviews, 95 that were smoker proxy interviews, and 102 that were originally self reports.

The reinterview study was restricted to those households with three or fewer smokers in order to reduce the difficulty of identifying the original respondent, since the names of individuals were not collected as part of the original survey. This restricted set comprised 95 percent of the households in the original study. Only one respondent was interviewed in any household at the reinterview. The person to be interviewed was identified by the original reporting status and by demographic information such as age, sex, and education. If there was any question as to whether the respondent was the person originally interviewed, the case was replaced. Similarly, if a respondent refused, no attempt was made to convert the refusal for the reinterview. Close to 100 interviews were completed in each of the three respondent groups. Because of the decision-rules, twice as many cases were randomly assigned as were ultimately thought to be needed.

The questionnaire used for the main study was done using CATI while the questionnaire used for the reinterview was done using hard copy. The questionnaire used for the reinterview contained all of the questions pertaining to cigarette-related information and a few demographic questions. The questions were worded exactly as they were worded in the main study. An introductory phrase was added to most questions which said "as of the date of the previous interview" to place the respondent in the context of the interview date.

Finally, respondents from each of the three groups were randomly assigned to one of two versions of the questionnaire. The only difference between the two versions was the wording of the categories for the income question. In version one, for example, a category reads "\$10,000 - 19,999 a year." In version two, the category reads "\$10,000 up to \$20,000." There was a special need in the study to test the subtle difference in wording. Both versions of the questionnaire appear as an attachment to this report.



#### IV. FINDINGS BASED ON THE REINTERVIEW DATA

This section examines the data from the random sample selected for reinterview, assessing the reliability of proxy responses provided in the main interview by comparing the responses on a follow-up survey of randomly selected smokers with the proxy responses obtained on the initial interview. In addition, because the sample includes reinterviews with some individuals who were interviewed themselves in the initial sample, the (test-retest) reliability of data is measured, and the reliability of responses by the type of the initial respondent can be compared. The degree to which individual data items are missing for the original survey and the reinterview survey is also examined. For the income question, the reinterview also tested two versions of the question that used slightly different wording.

Analysis of reinterview data included the variables measuring smoker behavior and cigarette characteristics, two household characteristics--number of smokers in household and household income--and smoker's age. The income variable was included because of interest in testing two versions of question wording.

The analysis examines first the degree to which reinterview responses match those of the initial survey and how this differs by type of initial respondent. The degree to which the reinterview respondent (always a self report) was able to provide data not reported by proxy respondents is then examined.

The percentage of mismatches varies across variables and across original respondent groups. Table IV.1 presents the percentage of responses that do not match, given that data was provided on both the original survey and the reinterview. Overall, the percentage of mismatches ranges from zero for whether the cigarette was filtered to 45 percent for income category. The percentages of mismatches for cigarette information ranges from zero for filtered to 32 percent for brand code. Several differences between groups are also seen. Except for household characteristics, the degree of mismatch is highest for cases where the original respondent was a non-smoker proxy.

TABLE IV.1

PERCENTAGE OF CASES WHERE REINTERVIEW RESPONSES DO NOT MATCH ORIGINAL RESPONSE,  
BY TYPE OF RESPONDENT TO ORIGINAL INTERVIEW

Variable	Percentage (%) Mismatch and Number of Cases							
	Original Respondent			Total Sample	$\chi^2$	Df	p =	Significant Contrast <sup>a</sup>
1 Self	2 Smoker Proxy	3 Nonsmoker Proxy						
1. # Smokers in HH Sample Size	11.1 99	28.4 95	14.4 97	17.9 291	11.07	2	0.004	1-2, 2-3
2. Brand Code Sample Size	25.5 98	34.4 93	35.7 90	32.0 281	3.04	2	0.218	None
3. Length of Cigarette (Regular/Long/Extra Long Sample Size)	7.1 98	19.6 92	23.3 90	15.4 280	9.94	2	0.007	1-2, 1-3
4. Filtered or Not Sample Size	0.0 98	0.0 95	0.0 90	0.0 283	NA	NA	NA	None
5. Pack Type (soft or hard) Sample Size	7.4 95	10.0 90	21.6 88	12.8 273	9.22	2	0.010	1-3, 2-3
6. Mentholated or Not Sample Size	4.1 98	9.5 95	6.1 87	6.8 280	2.22	2	0.330	None
7. Amount Smoked Per Day (Whether more than a pack) Sample Size	13.1 99	22.6 93	25.3 83	20.0 275	4.76	2	0.092	1-3
8. Annual HH Income (in \$10,000 intervals) Sample Size	35.1 77	64.4 73	36.7 68	45.4 218	6.33	2	0.042	1-2, 2-3
9. Age Within 2 Years Sample Size	8.2 98	7.5 93	15.2 92	10.3 283	93.68	2	0.159	NA

Includes only cases where a valid response (other than don't know) was provided on both surveys.

<sup>a</sup>Comparisons where the between group difference is significant at the 5 percent level. Contrast 1-2 is self vs. smoker proxy, 2-3 is smoker proxy vs. non-smoker proxy, 1-3 is self vs. non-smoker proxy.

Among the measures of smoking behavior and cigarette characteristics, the most notable differences across respondent groups are in the percent of mismatches on brand code, length of cigarette, pack type and amount smoked. For each of these variables the difference between the group with the highest mismatch and that with the lowest is 10 percentage points or more. However, only the differences for length and pack type were statistically significant at the 5 percent level. While the observed mismatch on these two variables was highest for the non-smoker proxies, the only large (significant at the 5 percent level) difference between the two proxy groups was for pack type.

The results must be evaluated in light of the degree of mismatch between the self-reports at the original and reinterview since that is as accurate as proxy responses can be expected to get. The degree of mismatch for the individuals who originally supplied data on themselves (self-respondents) is surprisingly high for some variables. For example, the self-mismatch for brand code is 25.5%, lower than the degree of mismatch for the two proxy groups (34.4 and 35.1 percent, respectively) but higher than what one might expect. Because brand code is perhaps the most essential cigarette characteristic collected, two factors will be examined to explain the degree of mismatch, namely:

- the difference in elapsed time between the original interview and the reinterview for the matches and mismatches
- the frequency of brand change cases as reported at the reinterview for matches and mismatches

Table IV.2 presents the mean number of days which elapsed between the original interview and the reinterview by type of case. The mean number of days which elapsed between the interview and the reinterview for the sample as a whole was 66.7; 68 for the original self reporter; 66.5 for the smoker proxy; and 65.7 for the nonsmoker proxy. The range and distributions for elapsed time were also similar about the same. The nonsmoker proxy had more mismatches on the whole and slightly less time elapsed between the original interview and reinterview. Similarly, those cases where the brand mismatched had the least number of elapsed days (63.5) between interviews. While there is some

**TABLE IV.2**

**NUMBER OF DAYS ELAPSING BETWEEN ORIGINAL INTERVIEW  
AND REINTERVIEW BY TYPE**

	<b>Mean</b>	<b>Median</b>
<b>Reinterview Sample as Whole</b>	<b>66.7</b>	<b>80</b>
<b>Original Self-Reporter</b>	<b>68.0</b>	<b>82</b>
<b>Original Smoker Proxy</b>	<b>66.5</b>	<b>78</b>
<b>Original Nonsmoker Proxy</b>	<b>65.7</b>	<b>80</b>
<b>Cases Brand Code Matched</b>	<b>68.4</b>	<b>82</b>
<b>Cases Brand Code Mismatched</b>	<b>63.5</b>	<b>45</b>

difference, one would expect more accurate and reliable data with the least amount of time elapsing between interviews. Because there is in fact less reliability with the least amount of elapsed time one might conclude that the amount of elapsed time between interviews does not explain the relatively high level of mismatch on brand code for the sample as a whole.

As part of the reinterview, respondents were asked how frequently they changed the brand of cigarette they smoked. Table IV.3 presents this data by type of original report and match or mismatch on brand code. While the number of cases in the most frequent categories are smaller, there is a pattern for the most frequent brand changers to have a greater percentage of mismatch than those who seldomly or never change their brand. This is as expected—if a person frequently changed their brand they would be less likely to recall what brand they were smoking two months or more before the interview. Also, while respondents were asked to report the usual brand they smoked, some respondents said they had no "usual" brand. In those cases, they were asked to report the brand they smoked most often and if that was not possible the brand they smoked closest to the interview. Individuals who had no usual brand may have reported accurately at the time of the interview but could not remember accurately at a later time. (Recall that respondents were asked to think back and report as of the date of the original interview.) Nonetheless, even among self respondents who say they never change brands, 20 percent gave a different brand at reinterview than they did in the initial interview.

Differences in a proxy's ability to report on the length of cigarette someone else smoked is somewhat understandable. This question provided three answer choices requiring a finer distinction of (1) regular or kings (2) long or deluxe and (3) extra long. While most of the others have two answers indicating the presence or absence of a characteristic. This information may be too refined for some of those proxy reporters reporting for others in a household. The degree to which this fact affects the use of the data for the 16.4 percent having a mismatch depends on how different the

TABLE IV.3

## FREQUENCY OF BRAND CHANGE BY TYPE OF REPORT AND BRAND CODE MATCH STATUS

	SELF		SMOKER PROXY		NONSMOKER PROXY	
	Match	Mismatch	Match	Mismatch	Match	Mismatch
<b>Frequently</b>						
Number	1	4	2	7	1	2
Percent	20	80	22	78	33	67
<b>Once in a While</b>						
Number	12	7	7	8	5	11
Percent	63	37	47	53	31	69
<b>Seldomly</b>						
Number	19	5	20	8	17	14
Percent	79	21	71	29	55	45
<b>Never</b>						
Number	41	10	32	11	34	13
Percent	80	20	75	25	72	28

cigarette characteristics (porosity, density, etc.) are when analyzed for these cases by length. In the worst cases, it affects fewer than 16.4 percent because some provided a UPC code which more accurately matches the data in any case.

The differences between the original responses and reinterview responses appear to differ randomly and are not systematically biased toward a particular response. This was assessed by crosstabulating the original survey responses for reinterview sample members with the reinterview response. The marginal distributions, presented in Table IV.4, are very similar for the two sources of data, and examination of the off-diagonal elements of the crosstabulation shows that the mismatches are very evenly distributed with the reinterview responses being equally likely to be shorter or longer than the original responses. Furthermore, this pattern occurs for all three groups of original respondent types. Thus, while the proportion of mismatches is higher for the both groups of proxy respondents, the overall distribution does not appear to have been affected by the differences.

The other cigarette characteristic for which the proportion of mismatches was significantly greater for proxies than for self-respondents was pack type (soft or hard). Again, the overall distribution is quite similar for the original survey response and the reinterview with the smokers themselves (Table IV.5), but the original respondents were slightly more likely to indicate soft pack than were the reinterview respondents. Examination of the original survey-reinterview crosstabulation for each of the three respondent groups separately shows that this pattern occurs for all three groups, including the group of original self respondents. While the proportion of mismatches is clearly lower for the self-respondents (7.4 percent) than for the two proxy groups (especially the non-smoker proxy group, at 21.6 percent), the pattern of a higher reported use of hard packs in the original interview than in the reinterview exists for all groups. Thus, the difference may be due more to the passage of time than an indication that non-smoker proxies at the original interview gave frequent incorrect responses.

TABLE IV.4

A COMPARISON OF THE DISTRIBUTIONS ON LENGTH OF CIGARETTE  
FOR THE ORIGINAL AND REINTERVIEW SAMPLE

Reinterview	Original Survey				Total
	Regular/ long	Long/ deluxe	Extra long	Don't know	
Regular/long	53%	7%	0%	2%	61.2%
Long/deluxe	7%	25%	1%	<1%	34.4%
Extra long	0%	1%	2%	<1%	3.8%
Don't know	<1%	0	0	<1%	0.7%
Total	60.5%	33.0%	3.1%	3.4%	100.0%



TABLE IV.5

COMPARISON OF ORIGINAL AND REINTERVIEW RESPONSES ON  
CIGARETTE PACKAGING FOR REINTERVIEW SAMPLE CASES

Reinterview	Original Survey Response			Total
	Soft Pack	Hard Pack	Don't know	
Soft Pack	63%	4%	2%	68.7%
Hard Pack	8%	19%	1%	28.5%
Don't know	1%	1%	1%	2.8%
Total	72.2%	23.7%	4.1%	100.0%

NOTE: Data are for 291 individuals.

For other measures, the percentage of mismatches was highest for income (45.4 percent) and lowest for age within 2 years (10.3 percent). For income and number of smokers in the household the highest degree of mismatch was for non-smoker proxies.

The mismatch on income category for smoker proxies is substantially higher at 64.6 percent than the 35.1 percent for the self-reports and 36.7 for nonsmoker proxies. This sizeable difference suggests that the three groups may differ on other personal characteristics which may be associated with knowledge of household income. For example, more self-reporters were women and survey experience indicates that more women answer the telephone. If this is the case, perhaps more male head of households were smoker proxies (although complete information is not available from the data set). Other studies indicate that more adult females answer "don't know" to household income questions and that when the answer is given it is often different from that reported by the male adult "head" of household. However, this is only one possible explanation for the high level of mismatch in the smoker proxy group for the income question.

Although the degree of mismatch was quite high for the income question the overall reliability was similar for the two versions of the income question (seen in Table IV.6) (43.1 percent for version 1 overall compared to 47.7 percent for version 2). Table IV.6 shows a larger discrepancy between the two versions within each of the three respondent groups than overall, however, ranging from a 15 percentage point difference when the original respondent was a smoker proxy, to a 9-11 percentage point difference for the other groups. The differences within subgroups are not large enough to be statistically significant due to small sample sizes.

Finally, the degree of mismatch on age was highest where the original respondent was a non-smoking proxy, but the difference between this group and the self respondents in percent mismatched is not large enough to be statistically significant at even the 10 percent level.

TABLE IV.6

NON-MATCH RESPONSES TO INCOME QUESTION BY VERSION, OVERALL  
AND BY TYPE OF ORIGINAL RESPONDENT

	Version 1	Version 2	Total	$\chi^2$	p =
<b>Percent Mismatch</b>					
<b>Overall</b> n =	43.1 109	47.7 109	45.4 218	0.46	0.496
<b>By Original Respondent</b>					
<b>Self Report</b> n =	30.8 39	39.5 38	35.1 77	0.64	0.424
<b>Smoker Proxy</b> n =	56.7 37	72.2 36	64.4 73	1.90	0.168
<b>Non-Smoker Proxy</b> n =	42.2 33	31.4 35	36.7 68	0.88	0.347

**NOTE:** The  $X^2$  statistic reported in the fifth column is for a test of whether the distributions of the responses for the three types of respondents differ by more than might be expected due to normal sampling variability, if the three samples had each been drawn from the same population. The p value in the last column gives the probability of observing a dispersion as large as that which is actually observed if the samples had been drawn from the same population.

Data on non-response are shown in Tables IV.7 and IV.8. The figures in these tables show the percentage of cases where responses were missing from:

- both the initial interview and the reinterview
- the reinterview only
- the original interview only

Data missing from both interviews indicate no change in the quality of data. If the original respondent was a proxy, data missing from the reinterview indicates that the proxy provided more information than the self-reporter at reinterview, while data missing from the original interview "only" indicate that the proxy provided less information. The amount of data that is missing is another indication of the relative quality of data provided by the three groups of original respondents.

The comparisons also indicate that non-smoker proxies were less likely than the other groups to provide data that the smoker would have provided as a self-reporter. Noteworthy differences are seen for several smoking measures: whether filtered or mentholated cigarettes are smoked, and amount smoked. Smaller differences are seen for length of cigarette and pack type. For other measures, the most noticeable result is the trivial difference on income. The difference on age of smoker is also small.

Table IV.8 presents a comparison of missing data by version of the income question. No differences are seen overall, and among the subgroups defined by original respondent, differences are seen only for non-smoking proxies, but these differences lead to no conclusions about whether one version is superior.

TABLE IV.7

MISSING DATA BY TYPE OF ORIGINAL RESPONDENT

	Self (n = 102)	Smoker Proxy (n = 95)	Non-Smoker Proxy (n = 98)	Total	$\chi^2$	p =
<u>Number of Smokers in Household</u>						
Data Missing From						
Both	0.0	0.0	0.0	0.0		
Reinterview	2.9	0.0	2.0	1.7		
Original	0.0	0.0	0.0	0.0	2.6	0.265
<u>Brand Code</u>						
Data Missing From						
Both	0.0	0.0	0.0	0.0		
Reinterview	2.9	0.0	1.0	1.4		
Original	0.0	0.0	0.0	0.0	3.31	0.192
<u>Length</u>						
Data Missing From						
Both	1.0	0.0	0.0	0.3		
Reinterview	2.9	0.0	2.0	1.7		
Original	0.0	3.2	6.1	3.1	10.80	0.095
<u>Filtered</u>						
Data Missing From						
Both	0.0	0.0	0.0	0.0		
Reinterview	2.9	0.0	1.0	1.4		
Original	1.0	0.0	7.1	2.7	14.39	0.006
<u>Pack Type</u>						
Data Missing From						
Both	2.9	0.0	1.0	1.4		
Reinterview	2.9	3.2	2.0	2.7		
Original	1.0	2.1	7.1	3.4	9.93	0.128

TABLE IV.7 (continued)

	Self (n = 102)	Smoker Proxy (n = 95)	Non-Smoker Proxy (n = 98)	Total	$\chi^2$	p =
<u>Mentholated</u>						
Data Missing From Both	1.0	0.0	1.0	0.7		
Reinterview	2.9	0.0	1.0	1.4		
Original	0.0	0.0	9.2	3.1	22.97	0.001
<u>Amount Smoked</u>						
Data Missing From Both	0.0	0.0	1.0	0.3		
Reinterview	2.9	0.0	2.0	1.7		
Original	0.0	2.1	12.2	4.8	22.62	0.001
<u>Income</u>						
Data Missing From Both	10.8	7.4	6.1	8.1		
Reinterview	7.8	7.4	9.8	9.8		
Original	9.8	8.4	10.2	8.1	6.00	0.424
<u>Age of Smoker</u>						
Data Missing From Both	0.0	0.0	0.0			
Reinterview	2.9	0.0	1.0			
Original	1.0	2.1	5.1		6.66	0.155

NOTE: The  $\chi^2$  statistic reported in the fifth column is for a test of whether the distributions of the responses for the three types of respondents differ by more than might be expected due to normal sampling variability, if the three samples had each been drawn from the same population. The p value in the last column gives the probability of observing a dispersion as large as that which is actually observed if the samples had been drawn from the same population.

TABLE IV.8

**NON-MATCH RESPONSES TO INCOME QUESTION BY VERSION, OVERALL  
AND BY TYPE OF ORIGINAL RESPONDENT**

	Version 1	Version 2	Total	$\chi^2$	p =
<b>Data Missing Overall</b>					
From Both	10.6	5.6	8.1	3.48	0.323
Reinterview	10.6	9.0	9.8		
Original	6.6	9.7	8.1		
n =	151	144	295		
<b>By Original Respondent Self</b>					
From Both	11.5	10.0	10.7	4.73	0.192
Reinterview	11.5	4.0	7.8		
Original	1.9	10.0	5.9		
n =	52	50	102		
<b>Smoker Proxy</b>					
From Both	8.1	6.5	7.4	0.705	0.872
Reinterview	6.1	8.7	7.4		
Original	10.2	6.5	8.4		
n =	49	46	95		
<b>Non-Smoker Proxy</b>					
From Both	12.0	0.0	6.1	6.42	0.093
Reinterview	14.0	14.6	14.3		
Original	8.0	12.5	10.2		
n =	50	48	98		

**NOTE:** The  $X^2$  statistic reported in the fifth column is for a test of whether the distributions of the responses for the three types of respondents differ by more than might be expected due to normal sampling variability, if the three samples had each been drawn from the same population. The p value in the last column gives the probability of observing a dispersion as large as that which is actually observed if the samples had been drawn from the same population.

## V. CONCLUSIONS

The rates of disagreement between the responses given by original proxy respondents and the responses subsequently elicited from the smokers themselves in the reinterview survey are fairly high for some of the questions; the overall range is from 0 to 64 percent. When compared to the percentage of mismatches among smokers who were interviewed initially and then reinterviewed, significantly higher rates of mismatch exist for five of the variables examined. One or both groups of proxy respondents had significantly higher rates of mismatches than the self respondents for two of the five cigarette characteristics, (length of cigarette, pack type), one of the two smoker characteristics (amount smoked), and both of the household characteristics (income, number of smokers in household).

These significant differences, however, appear to be reflect more on the design of the reinterview survey than on the quality of proxy responses at the original interview. That is, differences observed between the data supplied at reinterview and these supplied originally differ in larger part because of the passage of time (over two months on average) between the original survey and the reinterview, and to the change in the variable over time that may make it difficult to recall the appropriate response for an earlier point in time. Such problems of recall error are particularly likely for cigarette characteristics. For example, one-fourth of smokers report a different brand at reinterview than they reported themselves originally. Questions about cigarette characteristics may also have had ambiguous answers originally, further increasing the difficulty of recall. For example, many smokers may alternate between different lengths of cigarettes or pack type, depending upon what is readily available at the place of purchase. On the other hand, for variables that are likely to be fairly stable, such as whether the smoker buys filtered or unfiltered cigarettes, we observe no describable difference between original proxy and original self-respondents in the percentage of mismatches between the two interviews.

The higher rate of mismatches between original and reinterview responses for the groups with proxy respondents originally is therefore not surprising for the cigarette variables and does not necessarily mean that proxies *at the original interview* gave responses different from what the actual smoker would have



given *at that time*. If a smoker changed his type of cigarette or cigarette package frequently, he or she would clearly be better able than a proxy to remember the type of cigarette smoked or package purchased two or three months earlier. The proxy's response about smoking behavior at the time of interview may well have been nearly as accurate as the smoker's own response.

Mismatches between interviews for the household variables (income and number of smokers) were also higher when the original interview was with a smoker proxy, which is likely to be due simply to the consistency of the respondent rather than to systematic differences in the reliability of the response. Two different smokers in a household may well respond differently if asked about household income at any point in time. If one of these individuals were reasked about household income a few months later, the likelihood that the respondent will give an answer consistent with their own earlier response is greater than the likelihood that the respondent will give a response similar to the original response of the other smoker. However, there is no reason to believe *a priori* that the original respondent provided a more accurate estimate of household income than other smokers in the household would have given. The lower incomes typically reported at the reinterview than at the original interview with a smoker proxy suggests that these types of individuals may differ on a number of characteristics related to their knowledge of household incomes (e.g., original respondents may be more or less likely to be the head of household than those for whom a proxy provided the data initially). The difference between the original and reinterview responses for smoker proxy cases is due entirely to reinterview respondents indicating that there was only one smoker in the household. (By definition, smoker proxy cases were reported to have two or more smokers in the household on the original interview.)

In the eventual analysis of the effects of smoker and cigarette characteristics on the likelihood of a smoking related fire, more credible results will be obtained if the proxy responses were included than if they were excluded. While excluding proxy cases would eliminate any potential biases due to misreporting by proxies, these biases are likely to be relatively minor compared to the biases that would

be created by deleting these cases. If smokers who were the original respondents differ markedly from other smokers in these households, as the comparisons suggest that they do, deleting these cases from the analysis would yield a distorted sample of smokers and could lead to biased estimates of the relationship between smoker characteristics and smoking-related fires. Furthermore, the loss of proxy cases (one-fourth the sample if only nonsmoker proxy cases were deleted, one-half if both types of proxy cases were dropped) would substantially increase the variance of the estimates.

Two other arguments can also be made in favor of retaining the proxy cases. First, as indicated above, the differences between the original and reinterview responses exist only for some characteristics, and even for these the differences are not necessarily indicative of "errors" made by proxies in reporting for other smokers. Second, econometric studies suggest that the coefficients in linear regression models are "attenuated" (biased toward zero) when estimated on data with random errors in measurement. To the extent that the same effects occur in logit models, the bias in the estimates due to the measurement error would be to *understate* the effects of cigarette characteristics on fires. Thus, results which show a significant relationship would not be attributable to the measurement error and would be a conservative estimate of the true effects.

It is also recommended, however, that estimates be obtained with proxy cases removed, as a sensitivity test. If the cigarette characteristics continued to be significant predictors of the probability of a fire even when only one observation per household is used, this will provide support for the findings from the full model. Another sensitivity test that might be explored would be to select at random a single smoker from household with multiple smokers, to avoid any effects of inherent differences between original self-respondents and those for whom a proxy completed the interview.