

**SAS 2:**  
**Getting comfortable with your data**



**University of Guelph**

*Revised June 2011*

## Table of Contents

<b>SAS Availability</b> .....	<b>2</b>
<b>Data for SAS sessions</b> .....	<b>3</b>
<b>Review – Getting data into SAS</b> .....	<b>4</b>
Temporary SAS datasets .....	4
Permanent SAS datasets .....	4
<b>Statistical Refresher</b> .....	<b>5</b>
Types of Statistics .....	5
Types of Variables .....	5
Appropriate Statistics .....	6
<b>Frequency</b> .....	<b>7</b>
<b>Exercise 1</b> .....	<b>13</b>
<b>Mean, mode and median</b> .....	<b>13</b>
<b>Normality</b> .....	<b>16</b>
<b>Transformations</b> .....	<b>21</b>
<b>Exercise 2:</b> .....	<b>22</b>
<b>ODS – Output Delivery System</b> .....	Error! Bookmark not defined.
<b>SAS/GRAPH</b> .....	Error! Bookmark not defined.

# SAS Availability

Faculty, staff and students at the University of Guelph may access SAS three different ways:

## **1. Library computers**

On the library computers, SAS is installed on all machines.

## **2. Acquire a copy for your own computer**

If you are faculty, staff or a student at the University of Guelph, you may obtain the site-licensed standalone copy of SAS at a cost. However, it may only be used while you are employed or a registered student at the University of Guelph. To obtain a copy, go to the CCS Software Distribution Site ([www.uoguelph.ca/ccs/download](http://www.uoguelph.ca/ccs/download)).

# Goals of the workshop

This workshop builds on the skills and knowledge develop in "Getting your data into SAS". Participants are expected to have basic SAS skills and statistical knowledge. Specific goals of this workshop:

- To review reading data into SAS datasets
- To learn how to determine whether your data comes from a Normal distribution
- How do we transform data if it is needed
- Plotting your data using SAS/GRAPH procedures

# Data for SAS sessions

## **Dataset: Canadian Tobacco Use Monitoring Survey 2010 – Person File**

This survey tracks changes in smoking status, especially for populations most at risk such as the 15- to 24-year-olds. It allows Health Canada to estimate smoking prevalence for the 15- to 24-year-old and the 25-and-older groups by province and by gender on a semi-annual basis.

The sample data used for this series of SAS workshops only includes respondents from the province of Quebec and only 14 of a possible 202 variables are being used.

To view the data, open the Excel spreadsheet entitled CTUMS\_2010.xls

Variable Name	Label for Variable
<b>PUMFID</b>	Individual identification number
<b>PROV</b>	Province of the respondent
<b>DVURBAN</b>	Characteristic of the community
<b>HHSIZE</b>	Number of people in the household
<b>HS_Q20</b>	Number of people that smoke inside the house
<b>DVAGE</b>	Age of respondent
<b>SEX</b>	Respondent's sex
<b>DVMARST</b>	Grouped marital status of respondent
<b>PS_Q30</b>	Age smoked first cigarette
<b>PS_Q40</b>	Age begin smoking cigarettes daily
<b>WP_Q10A</b>	Number of cigarettes smoked – Monday
<b>WP_Q10B</b>	Number of cigarettes smoked – Tuesday
<b>WP_Q10C</b>	Number of cigarettes smoked – Wednesday
<b>WP_Q10D</b>	Number of cigarettes smoked – Thursday
<b>WP_Q10E</b>	Number of cigarettes smoked – Friday
<b>WP_Q10F</b>	Number of cigarettes smoked – Saturday
<b>WP_Q10G</b>	Number of cigarettes smoked – Sunday
<b>SC_Q100</b>	What was the main reason you began to smoke again?
<b>WTPP</b>	Person weight (survey weight variable)

# Review – Getting data into SAS

## Temporary SAS datasets

When you create/save a dataset in SAS, whether it by using an infile statement, cards statement or the import procedure in SAS, by default SAS places it in the Work library. You will not see a physical SAS dataset file on your computer when you use the Work library – in other words the SAS dataset that is created is a Temporary SAS dataset.

## Permanent SAS datasets

There may be situations when you may need to read from a permanent SAS dataset or you may need to create a physical SAS dataset file. This will require the use of the SET and LIBNAME statement. The SET statement refers to the filename of the SAS data set already created or referred to in this SAS session. LIBNAME refers to a location where you would like to save your SAS data set.

```
libname sasdata "C:\Users\edwardsm\Documents\Michelle_Docs\Workshops";
```

```
Data sasdata.ctums2;  
set ctums;  
Run;
```

The dataset name will now need to use its full name which includes its library name – sasdata in this case.

The Set statement tells SAS to use the dataset we've already created and stored in the Work library of SAS

The libname statement lets SAS know WHERE you want to save the SAS dataset. In this example a directory called Workshop in a Michelle\_Docs directory

Now we will have 2 SAS datasets – one in the Work Library called ctums and a second called ctums2 in the SASDATA library which is located in the C:\Users\edwardsm\Documents\Michelle\_Docs\Workshops directory. If you look in the specified directory you should now see a file called ctums2.sas7bdat. This is referred to as a permanent SAS dataset – there is now a physical file that we can see and send to colleagues.

When someone sends you a \*.sas7bdat file how do you read it?

1. You will need to create a Library – by using the libname statement. The name of the Library can be anything you'd like it to be. The name of the library is only used to create a location on your computer and to refer to it in your program.
2. By using a Data step – create either a local copy (Work library) or save it as a different name.

*Example:*

I have just received the file ctums2.sas7bdat in my email. I will save it in a new directory called C:\Users\edwardsm\Documents\Michelle\_Docs\Research

```
libname newdata "C:\Users\edwardsm\Documents\Michelle_Docs\Research";
```

```
Data newctums;  
set newdata.ctums2;  
Run;
```

My permanent SAS dataset ctums2.sas7bdat is located in the Research directory on my computer. I create the library called newdata which refers to that directory.

I save a new copy called newctums in my Work library.

# Statistical Refresher

## Types of Statistics

Two broad types of statistics exist which are descriptive and inferential. Descriptive statistics describe the basic characteristics of the data in a study. Usually generated through an Exploratory Data Analysis (EDA), they provide simple numerical and graphical summaries about the sample and the measures. Inferential statistics allow you to make conclusions regarding the data i.e. significant differences, relationships between variables, etc.

Here are some examples of descriptive and inferential statistics:

Descriptives	Inferential
<ul style="list-style-type: none"><li>• Frequencies</li><li>• Means</li><li>• Standard Deviations</li><li>• Ranges</li><li>• Medians</li><li>• Modes</li></ul>	<ul style="list-style-type: none"><li>• T-tests</li><li>• Chi-squares</li><li>• ANOVA</li><li>• Friedman</li></ul>

Which test to perform on your data largely depends on a number of factors including:

1. What type of data you are working with?
2. Are your samples related or independent?
3. How many samples are you comparing?

## Types of Variables

Variable types can be distinguished by various levels of measurement which are Nominal, Ordinal, Interval or Ratio.

### Nominal

Have data values that identify group membership. The only comparisons that can be made between variable values are equality and inequality. Examples of nominal measurement include gender, race, religious affiliation, telephone area codes or country of residence.

### Ordinal

Have data values arranged in a rank ordering with an unknown difference between adjacent values. Comparisons of greater and less can be made and in addition to equality and inequality. Examples include: results of a horse race, level of education or satisfaction/attitude questions.

## **Interval**

Are measured on a scale such that a one-unit change represents the same difference throughout the scale. These variables do not have true zero points. Examples include: temperature in the Celsius or Fahrenheit scale, year date in a calendar or IQ test results.

## **Ratio**

Have the same properties as interval variables plus the additional property of a true zero. Examples include: temperature measured in Kelvins, most physical quantities such as mass, length or energy, age, length of residence in a given place.

Interval and Ratio will be considered identical thus yielding three types of measurement scales.

## **Appropriate Statistics**

For each type of variable a particular measure of central tendency is most appropriate. By central tendency we mean one value that most effectively summarizes a variable's complete distribution.

Measurement Scale	Measure of Central Tendency
Nominal	Mode – value that appears the most often in distribution.
Ordinal	Median – Value that divides the ordered distribution of responses into two equal size groups. (the value of the 50 <sup>th</sup> percentile)
Interval / Scale	Mean – The arithmetic average of a distribution.

# Frequency

How many males and females are in this dataset?

**SAS Code:**

```
Proc freq data=newctums;  
  tables sex;  
Run;
```

**Tables** statement – list the variables you would like to see a frequency chart created.

Good SAS code writing etiquette – use a **data=** option in your **Proc** statement – this ensures that SAS uses the correct dataset in your analysis.

Ensure that the procedure is closed with a **Run;**

**SAS Output:**

The FREQ Procedure

SEX

SEX	Frequency	Percent	Cumulative Frequency	Cumulative Percent
1	432	46.96	432	46.96
2	488	53.04	920	100.00



Based on the results we see that in our sample of CTUMS2010 47% of the sample are males and 53% are females.

Let's look at age\_group – what is the frequency distribution for the age\_group variable?

**SAS Code:**

```
Proc freq data=newctums;  
  tables agegroup;  
Run;
```

**SAS Output:**

The FREQ Procedure

agegroup	Frequency	Percent	Cumulative Frequency	Cumulative Percent
15-24 years	95	53.67	95	53.67
25-34 years	15	8.47	110	62.15
35-44 years	20	11.30	130	73.45
45-54 years	18	10.17	148	83.62
55-64 years	24	13.56	172	97.18
65-74 years	4	2.26	176	99.44
75-84 years	1	0.56	177	100.00

Of the sample, 53.67% are between the ages of 15 and 24 years.

Now let's put the two tables together and create a Cross-tabulation to show us the age distribution of the 2 genders. To accomplish this task we will list the 2 variables of interest in the **Tables** statement and place an '\*' between the two to let SAS know that we want a crosstab.

**SAS Code:**

```
Proc freq data=newctums;
  tables agegroup*sex;
Run;
```

Note: The order of the variables in your **Tables** statement determine the structure of the table.

Row variable \* column variable

**SAS Output:**

The FREQ Procedure  
Table of agegroup by SEX

agegroup	SEX(SEX)		Total
	Male	Female	
15-24 years	47	48	95
	26.55	27.12	53.67
	49.47	50.53	
	52.22	55.17	
25-34 years	9	6	15
	5.08	3.39	8.47
	60.00	40.00	
	10.00	6.90	
35-44 years	7	13	20
	3.95	7.34	11.30
	35.00	65.00	
	7.78	14.94	

There are 6 females between the ages of 25 and 34 yrs.

3.39 % of this sample are females between the ages of 25 and 34 yrs.

Of all the individuals between 25 and 34 years, 40 % are female.

Of all the females, 6.90 % are between the ages of 25 and 34 yrs.

45-54 years	10 5.65 55.56 11.11	8 4.52 44.44 9.20	18 10.17
55-64 years	13 7.34 54.17 14.44	11 6.21 45.83 12.64	24 13.56
65-74 years	3 1.69 75.00 3.33	1 0.56 25.00 1.15	4 2.26
75-84 years	1 0.56 100.00 1.11	0 0.00 0.00 0.00	1 0.56
Total	90 50.85	87 49.15	177 100.00

There are a total of 24 individuals between the ages of 55 and 64 yrs, which makes up 13.56% of the sample.

### Chi-square tests

If we want to test whether a relationship exists between 2 categorical variables, a Chi-square test is one option. To conduct a Chi-square test in SAS we will add an option to the above coding.

```
Proc freq data=newctums;
  tables agegroup*sex /chisq;
Run;
```

The FREQ Procedure

Table of agegroup by SEX

agegroup	SEX(SEX)		Total
	Male	Female	
55-64 years	13	11	24
	7.34	6.21	13.56
	54.17	45.83	
	14.44	12.64	
65-74 years	3	1	4
	1.69	0.56	2.26
	75.00	25.00	
	3.33	1.15	
75-84 years	1	0	1
	0.56	0.00	0.56
	100.00	0.00	
	1.11	0.00	
Total	90	87	177
	50.85	49.15	100.00

Statistics for Table of agegroup by SEX

Statistic	DF	Value	Prob
Chi-Square	6	4.7499	0.5763
Likelihood Ratio Chi-Square	6	5.2141	0.5167
Mantel-Haenszel Chi-Square	1	0.6105	0.4346
Phi Coefficient		0.1638	
Contingency Coefficient		0.1617	
Cramer's V		0.1638	

WARNING: 29% of the cells have expected counts less than 5. Chi-Square may not be a valid test.

Sample Size = 177

Note the Warning!!!  
With cells of 5 or less the Chi-square test may not be a valid test. Think about ways to recode your data ie. Create new groupings or re-examine your choice of statistical test

```
Proc freq data=newctums;
  tables dvmarst*sex /chisq;
Run;
```

Table of DVMARST by SEX

DVMARST(DVMARST)	SEX(SEX)		Total
	Male	Female	
Frequency			
Percent			
Row Pct			
Col Pct			
Common-law/Married	31 17.61 51.67 34.83	29 16.48 48.33 33.33	60 34.09
Widow/Divorced/Separated	5 2.84 38.46 5.62	8 4.55 61.54 9.20	13 7.39
Single	53 30.11 51.46 59.55	50 28.41 48.54 57.47	103 58.52
Total	89 50.57	87 49.43	176 100.00

Statistics for Table of DVMARST by SEX

Statistic	DF	Value	Prob
Chi-Square	2	0.8237	0.6624
Likelihood Ratio Chi-Square	2	0.8299	0.6604
Mantel-Haenszel Chi-Square	1	0.0017	0.9671
Phi Coefficient		0.0684	
Contingency Coefficient		0.0683	
Cramer's V		0.0684	

The non-significant chi-square suggests that there is no association between marital status and sex in this sample.

Sample Size = 176

# Exercise 1

**Exercise:** Is there an association between sex (**sex**) and whether the individual lived in an urban or rural area (**dvurban**)?

## Mean, mode and median

We've seen earlier that to determine the measures of central tendency – mean, median and mode are the best statistics. All three measures are available in **Proc UNIVARIATE** (see next section), but **Proc MEANS** also offers the mean and median.

### Mean

We will calculate the Mean for total number of cigarettes smoked in a week (**totcig**)

#### SAS Code:

```
Proc means data=newctums;  
  var totcig;  
Run;
```

#### SAS Output:

The MEANS Procedure

Analysis Variable : totcig

N	Mean	Std Dev	Minimum	Maximum
176	108.1079545	137.3571143	0	691.0000000

#### How do you get the standard error in this output?

## Median

We will calculate the Median of SC\_Q100 (What was the main reason you began to smoke again?)

### SAS Code:

```
Proc means data=newctums median;  
  var sc_q100;  
Run;
```

### SAS Output:

```
                The MEANS Procedure  
  
Analysis Variable : SC_Q100 SC_Q100  
  
                Median  
-----  
                4.0000000  
-----
```

## Mode

**Proc UNIVARIATE** is the only procedure that will calculate the Mode of a variable. The mode is defined as the value with the most observations – so you can accomplish with **Proc FREQ** as well.

### SAS Code (Proc UNIVARIATE):

```
Proc univariate data=newctums;  
  var sc_q100;  
Run;
```

**SAS Output:**

The UNIVARIATE Procedure  
Variable: SC\_Q100 (SC\_Q100)

Moments

N	72	Sum Weights	72
Mean	4.9027778	Sum Observations	353
Std Deviation	3.00778203	Variance	9.04675274
Skewness	0.53048708	Kurtosis	-1.3134182
Uncorrected SS	2373	Corrected SS	642.319444
Coeff Variation	61.3485287	Std Error Mean	0.35447051

Basic Statistical Measures

Location		Variability	
Mean	4.902778	Std Deviation	3.00778
Median	4.000000	Variance	9.04675
Mode	2.000000	Range	9.00000
		Interquartile Range	6.00000

Tests for Location: Mu0=0

Test	-Statistic-	-----p Value-----
Student's t	t 13.83127	Pr >  t  <.0001
Sign	M 36	Pr >=  M  <.0001
Signed Rank	S 1314	Pr >=  S  <.0001



## Normality

One assumption we see many tests demanding is that the dependent variable (if it is a continuous variable) comes from a normal distribution. The “bell curve” is often used as an example of a normal distribution. How do we go about testing this assumption in SAS?

One powerful **Procedure** to provide you with descriptive statistics along with normality tests in **UNIVARIATE**. To view all the syntax and options for this procedure please visit:

[http://support.sas.com/onlinedoc/913/getDoc/en/procstat.hlp/univariate\\_index.htm#idxuni0002](http://support.sas.com/onlinedoc/913/getDoc/en/procstat.hlp/univariate_index.htm#idxuni0002)

Let’s take a look at our “totcig” variable – this is the variable we created in SAS 1 that was the sum of cigarettes smoked during the week.

Let’s look at the default output for **Proc UNIVARIATE** first:

### SAS Code:

```
Proc univariate data=newctums;  
  var totcig;  
Run;
```

### SAS Output:

The UNIVARIATE Procedure  
Variable: totcig

Moments

N	176	Sum Weights	176
Mean	108.107955	Sum Observations	19027
Std Deviation	137.357114	Variance	18866.9769
Skewness	2.92438367	Kurtosis	9.10766212
Uncorrected SS	5358691	Corrected SS	3301720.95
Coeff Variation	127.055511	Std Error Mean	10.3536821

The moments provide you with a lot of descriptive statistics – from the Mean to Stderr to Kurtosis.

Basic Statistical Measures

Location		Variability	
Mean	108.1080	Std Deviation	137.35711
Median	70.0000	Variance	18867
Mode	84.0000	Range	691.00000

Interquartile Range 105.00000

Tests for Location: Mu0=0

Test	-Statistic-	-----p Value-----
Student's t	t 10.4415	Pr >  t  <.0001
Sign	M 85	Pr >=  M  <.0001
Signed Rank	S 7267.5	Pr >=  S  <.0001

Three test for location test the hypothesis that the sample mean is equal to 0

Quantiles (Definition 5)

Quantile	Estimate
100% Max	691
99%	691
95%	497
90%	185
75% Q3	140
50% Median	70
25% Q1	35
10%	3
5%	1
1%	0
0% Min	0

Extreme Observations

----Lowest----		----Highest---	
Value	Obs	Value	Obs
0	174	608	111
0	155	691	66
0	118	691	77
0	85	691	107
0	71	691	136

A listing of extreme observations is also provided.

By reviewing this output, the researcher now has a very good picture of their data. However, let's add two options, one to provide plots (**plot**) and a second to provide Normality test statistics (**normal**).

**SAS Code:**

```
Proc univariate data=newctums normal plot;  
  var totcig;  
Run;
```

**SAS Output:**

The UNIVARIATE Procedure  
Variable: totcig

Moments

N	176	Sum Weights	176
Mean	108.107955	Sum Observations	19027
Std Deviation	137.357114	Variance	18866.9769
Skewness	2.92438367	Kurtosis	9.10766212
Uncorrected SS	5358691	Corrected SS	3301720.95
Coeff Variation	127.055511	Std Error Mean	10.3536821

The first sections of the output match what we have seen above. With one exception:

Tests for Normality:

Basic Statistical Measures

Location		Variability	
Mean	108.1080	Std Deviation	137.35711
Median	70.0000	Variance	18867
Mode	84.0000	Range	691.00000
		Interquartile Range	105.00000

Tests for Location: Mu0=0

Test	-Statistic-	-----p Value-----
Student's t	t 10.4415	Pr >  t  <.0001
Sign	M 85	Pr >=  M  <.0001
Signed Rank	S 7267.5	Pr >=  S  <.0001

### Tests for Normality

Test	Statistic	p Value
Shapiro-Wilk	W 0.641166	Pr < W <0.0001
Kolmogorov-Smirnov	D 0.215624	Pr > D <0.0100
Cramer-von Mises	W-Sq 2.924821	Pr > W-Sq <0.0050
Anderson-Darling	A-Sq 17.57121	Pr > A-Sq <0.0050

### Quantiles (Definition 5)

Quantile	Estimate
100% Max	691
99%	691
95%	497
90%	185
75% Q3	140
50% Median	70
25% Q1	35

### Extreme Observations

Lowest		Highest	
Value	Obs	Value	Obs
0	174	608	111
0	155	691	66
0	118	691	77
0	85	691	107
0	71	691	136

### Tests for Normality:

The Null hypothesis for this test is – your dependent variable comes from a Normal distribution. So if the p-value < 0.05 then you reject the Null hypothesis and your dependent variable is not normal.

Four tests are provided.

Shapiro-Wilk should only be used if the sample size is less than or equal to 2000.



# Transformations

What happens when your data is not normally distributed as the above example? The first step is to try a transformation. Attempt to transform your data so it follows a normal distribution.

## Common Transformations:

1. Square root – commonly used for counts (**sqrt**)
2. Arcsin - commonly used for proportions (**arsin**)
3. Logarithmic transformation – log base 10 or natural log (**log10** or **log**)

For our data we will try the logarithmic transformation. Any data manipulations must be conducted in a **Data** step.

### SAS Code:

```
Data t_ctums;  
  set newctums;  
  log_totcig = log(1+totcig);  
Run;
```

**Data** – saving the new dataset under a new name “t\_ctums” – think of this as the File-SaveAs – option in Excel

**Set** - telling SAS to use the dataset already in the program – same function as if we were reading in a Permanent SAS dataset.

**Log\_totcig** - the name of a new variable we are going to create

**Log(totcig)** – totcig is the variable in our dataset that we want to transform – and log is the transformation function we want to use.

### Sample output:

Obs	totcig	log_ totcig
1	2	1.09861
2	105	4.66344
3	140	4.94876
4	84	4.44265

## Exercise 2:

- a) Rerun the Normality test to determine whether the new variable you created, log\_totcig, is indeed normal.
- b) What variable will you use to run your statistical tests? totcig or log\_totcig?
- c) How will you report your results?