

# Handout 1 SAS: Introduction to data

## Getting Started

Load the data set of 20,000 observations

```
In [1]: filename cdc url 'http://www.openintro.org/stat/data/cdc.csv';

proc import datafile=cdc
            out=work.cdc
            dbms=csv
            replace;
            getnames=yes;
run;
```

SAS Connection established. Subprocess id is 18327

Out[1]:

```
34 ods listing close;ods html5 (id=saspy_
internal) file=stdout options(bitmap_mode
='inline') device=svg; ods graphics on /
34 ! outputfmt=png;
NOTE: Writing HTML5(SASPY_INTERNAL) Body fi
le: STDOUT
35
36 filename cdc url 'http://www.openintr
o.org/stat/data/cdc.csv';
37
38 proc import datafile=cdc
39             out=work.cdc
40             dbms=csv
41             replace;
42             getnames=yes;
43 run;
NOTE: Unable to open SASUSER.PROFILE. WORK.
PROFILE will be opened instead.
NOTE: All profile changes will be lost at t
he end of the session.
44 /*****
*****
45 *   PRODUCT:    SAS
46 *   VERSION:    9.4
47 *   CREATOR:    External File Interfac
e
48 *   DATE:       07JUL18
49 *   DESC:       Generated SAS Datastep
Code
50 *   TEMPLATE SOURCE: (None Specifie
```

```

d.)
51      *****
*****/
52      data WORK.CDC      ;
53      %let _EFIERR_ = 0; /* set the ERRO
R detection macro variable */
54      infile CDC delimiter = ',' MISSOVE
R DSD lrecl=32767 firstobs=2 ;
55      informat genhlth $9. ;
56      informat exerany best32. ;
57      informat hlthplan best32. ;
58      informat smoke100 best32. ;
59      informat height best32. ;
60      informat weight best32. ;
61      informat wt desire best32. ;
62      informat age best32. ;
63      informat gender $1. ;
64      format genhlth $9. ;
65      format exerany best12. ;
66      format hlthplan best12. ;
67      format smoke100 best12. ;
68      format height best12. ;
69      format weight best12. ;
70      format wt desire best12. ;
71      format age best12. ;
72      format gender $1. ;
73      input
74          genhlth $
75          exerany
76          hlthplan
77          smoke100
78          height
79          weight
80          wt desire
81          age
82          gender $
83      ;
84      if _ERROR_ then call symputx('_EFI
ERR_',1); /* set ERROR detection macro var
iable */
85      run;
NOTE: The infile CDC is:
      Filename=http://www.openintro.org/sta
t/data/cdc.csv,
      Local Host Name=localhost.localdomai
n,
      Local Host IP addr=::1,
      Service Hostname Name=www.openintro.o
rg,
      Service IP addr=192.185.65.127,
      Service Name=httpd,Service Portno=80,
      Lrecl=32767,Recfm=Variable

NOTE: 20000 records were read from the infi
le CDC.
      The minimum record length was 24.
      The maximum record length was 31.

```

*NOTE: The data set WORK.CDC has 20000 observations and 9 variables.*

*NOTE: DATA statement used (Total process time):*

real time	1.53 seconds
cpu time	0.09 seconds

20000 rows created in WORK.CDC from CDC.

*NOTE: WORK.CDC data set was successfully created.*

*NOTE: The data set WORK.CDC has 20000 observations and 9 variables.*

*NOTE: PROCEDURE IMPORT used (Total process time):*

real time	2.86 seconds
cpu time	0.20 seconds

86

87 ods html5 (id=saspy\_internal) close;ods listing;

88

View the names of the variables

In [2]: 

```
proc contents data=work.cdc short;
run;
```

Out[2]: **The SAS System**

**The CONTENTS Procedure**

**Alphabetic List of Variables for WORK.CDC**

age exerany gender genhlth height hlthplan  
smoke100 weight wt desire

## Exercise 1

How many cases are there in this data set? How many variables? For each variable, identify its data type (for example, categorical, numeric).

There are 9 variables.

Data Source: US CDC website.

- genhlth: A categorical vector indicating general health, with categories excellent, very good, good, fair, and poor.
- exerany: A categorical vector, 1 if the respondent exercised in the past month and 0 otherwise.
- hlthplan: A categorical vector, 1 if the respondent has some form of health coverage and 0 otherwise.
- smoke100: A categorical vector, 1 if the respondent has smoked at least 100 cigarettes in their entire life and 0 otherwise.
- height: A numerical vector, respondent's height in inches.
- weight: A numerical vector, respondent's weight in pounds.
- wtdesired: A numerical vector, respondent's desired weight in pounds.
- age: A numerical vector, respondent's age in years.

Look at the first 10 rows of our data

```
In [3]: proc print data=work.cdc (obs=10);
run;
```

Out[3]: **The SAS System**

Obs	genhlth	exerany	hlthplan	smoke100
1	good	0	1	0
2	good	0	1	1
3	good	1	1	1
4	good	1	1	0
5	very good	0	1	0
6	very good	1	1	0
7	very good	1	1	0
8	very good	0	1	0
9	good	0	1	1
10	good	1	1	0

## Summaries and Tables

Look at numerical summary, such as mean, variance, standard deviation, minimum, maximum, and extreme observations

Summary statistic for weight

```
In [4]: proc univariate data=work.cdc;
        var weight;
        run;
```

Out[4]: **The SAS System**

**The UNIVARIATE Procedure**  
Variable: weight

<b>Moments</b>			
<b>N</b>	20000	<b>Sum Weights</b>	20000
<b>Mean</b>	169.68295	<b>Sum Observations</b>	3393659
<b>Std Deviation</b>	40.08097	<b>Variance</b>	1606.4797
<b>Skewness</b>	0.95572799	<b>Kurtosis</b>	1.99
<b>Uncorrected SS</b>	607974147	<b>Corrected SS</b>	321247.9
<b>Coeff Variation</b>	23.6210945	<b>Std Error Mean</b>	0.2811

<b>Basic Statistical Measures</b>			
<b>Location</b>		<b>Variability</b>	
<b>Mean</b>	169.6830	<b>Std Deviation</b>	40.08097
<b>Median</b>	165.0000	<b>Variance</b>	1606
<b>Mode</b>	160.0000	<b>Range</b>	432.00000

Basic Statistical Measures			
Location		Variability	
		Interquartile Range	50.00000

Tests for Location: $\mu_0=0$				
Test	Statistic		p Value	
Student's t	t	598.7079	Pr >  t	<.0001
Sign	M	10000	Pr >=  M	<.0001
Signed Rank	S	1.0001E8	Pr >=  S	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	500
99%	290
95%	240
90%	220
75% Q3	190
50% Median	165
25% Q1	140
10%	124
5%	115
1%	100
0% Min	68

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
68	18743	400	2944
70	16531	400	19319
78	18065	405	15720
78	11299	495	4445
79	7614	500	1995

The sample frequency distribution for smoke100

```
In [5]: proc freq data=work.cdc;
         tables Smoke100;
         run;
```

Out[5]: **The SAS System**

#### The FREQ Procedure

smoke100	Frequency	Percent	Cumulative Frequency
0	10559	52.80	10559
1	9441	47.21	20000

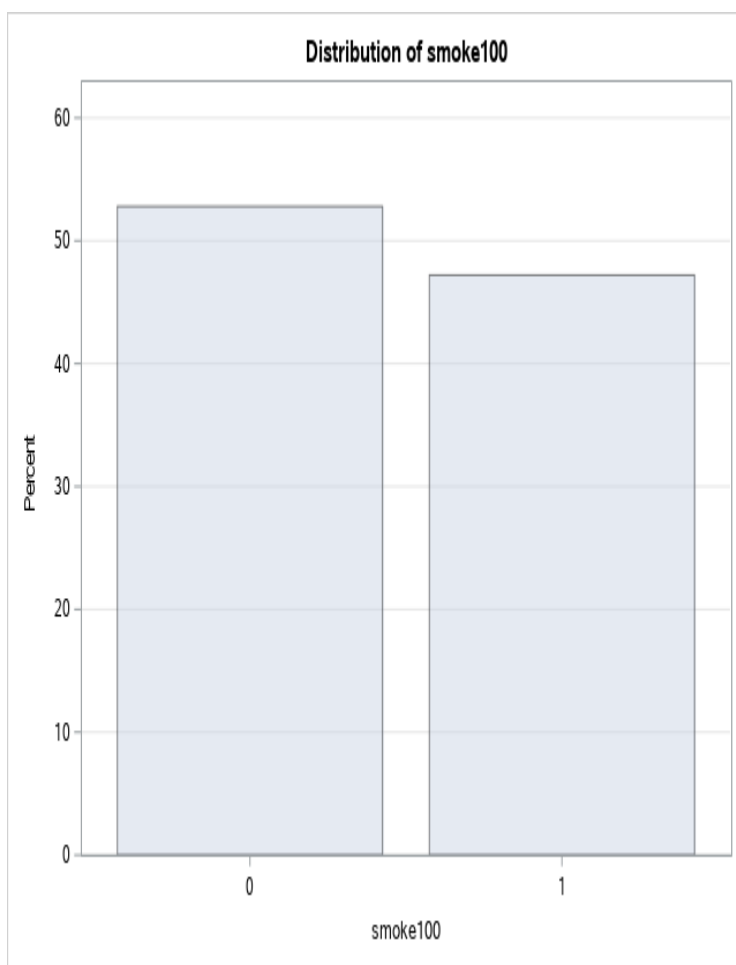
Graph the sample frequency distribution for smoke100

```
In [6]: proc freq data=work.cdc;
         tables Smoke100 / plots=freqplot(scale=percent);
         run;
```

Out[6]: **The SAS System**

## The FREQ Procedure

smoke100	Frequency	Percent	Cumulative Frequency
0	10559	52.80	10559
1	9441	47.21	20000



## Exercise 2

Create a numerical summary for height and age, and compute the interquartile range for each. Compute the relative frequency distribution for gender and exerany. How many males are in the sample? What proportion of the sample reports being in excellent health?

There are 9569 males in the sample. A total of 4657 out of 20000 reported being in excellent health. This is 23.29% of the sample.



## Numerical summary for height and age

```
In [7]: proc univariate data=work.cdc;
        var height;
        var age;
        run;
```

Out[7]:

## The SAS System

## The UNIVARIATE Procedure

Variable: height

Moments			
<b>N</b>	20000	<b>Sum Weights</b>	20000
<b>Mean</b>	67.1829	<b>Sum Observations</b>	1343658
<b>Std Deviation</b>	4.12595429	<b>Variance</b>	17.02350
<b>Skewness</b>	0.1036124	<b>Kurtosis</b>	-0.37095
<b>Uncorrected SS</b>	90611294	<b>Corrected SS</b>	340400
<b>Coeff Variation</b>	6.14137569	<b>Std Error Mean</b>	0.02054

Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	67.18290	<b>Std Deviation</b>	4.12595
<b>Median</b>	67.00000	<b>Variance</b>	17.02350
<b>Mode</b>	66.00000	<b>Range</b>	45.00000
		<b>Interquartile Range</b>	6.00000

## Tests for Location: Mu0=0

Tests for Location	Statistic	Ma0=0	p Value
--------------------	-----------	-------	---------

Test	Statistic		p Value	
Student's t	t	2302.763	Pr >  t	<.0001
Sign	M	10000	Pr >=  M	<.0001
Signed Rank	S	1.0001E8	Pr >=  S	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	93
99%	76
95%	74
90%	73
75% Q3	70
50% Median	67
25% Q1	64
10%	62
5%	61
1%	59
0% Min	48

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs

<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
48	15465	82	10160
48	5412	82	10322
49	8871	83	3691
50	3905	84	18817
51	11948	93	17534

### The SAS System

#### The UNIVARIATE Procedure

Variable: age

<b>Moments</b>			
<b>N</b>	20000	<b>Sum Weights</b>	20000
<b>Mean</b>	45.06825	<b>Sum Observations</b>	901365
<b>Std Deviation</b>	17.1926895	<b>Variance</b>	295.616
<b>Skewness</b>	0.45170032	<b>Kurtosis</b>	-0.6405
<b>Uncorrected SS</b>	46534419	<b>Corrected SS</b>	5911.5
<b>Coeff Variation</b>	38.1481186	<b>Std Error Mean</b>	0.12

<b>Basic Statistical Measures</b>	
<b>Location</b>	<b>Variability</b>

Basic Statistical Measures			
Location		Variability	
Mean	45.06825	Std Deviation	17.19269
Median	43.00000	Variance	295.58857
Mode	40.00000	Range	81.00000
		Interquartile Range	26.00000

Tests for Location: Mu0=0				
Test	Statistic		p Value	
Student's t	t	370.7165	Pr >  t	<.0001
Sign	M	10000	Pr >=  M	<.0001
Signed Rank	S	1.0001E8	Pr >=  S	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	99
99%	84
95%	76
90%	71
75% Q3	57
50% Median	43
25% Q1	31

**Quantiles (Definition 5)**

Level	Quantile
10%	24
5%	21
1%	18
0% Min	18

**Extreme Observations**

Lowest		Highest	
Value	Obs	Value	Obs
18	19885	95	16084
18	19860	96	17051
18	19832	97	10350
18	19706	99	900
18	19622	99	6710

Compute the interquartile range for height and age

Interquartile range for height: between 48 and 93 with a range of 45  
 Interquartile range for age: between 18 and 99 with a range of 81

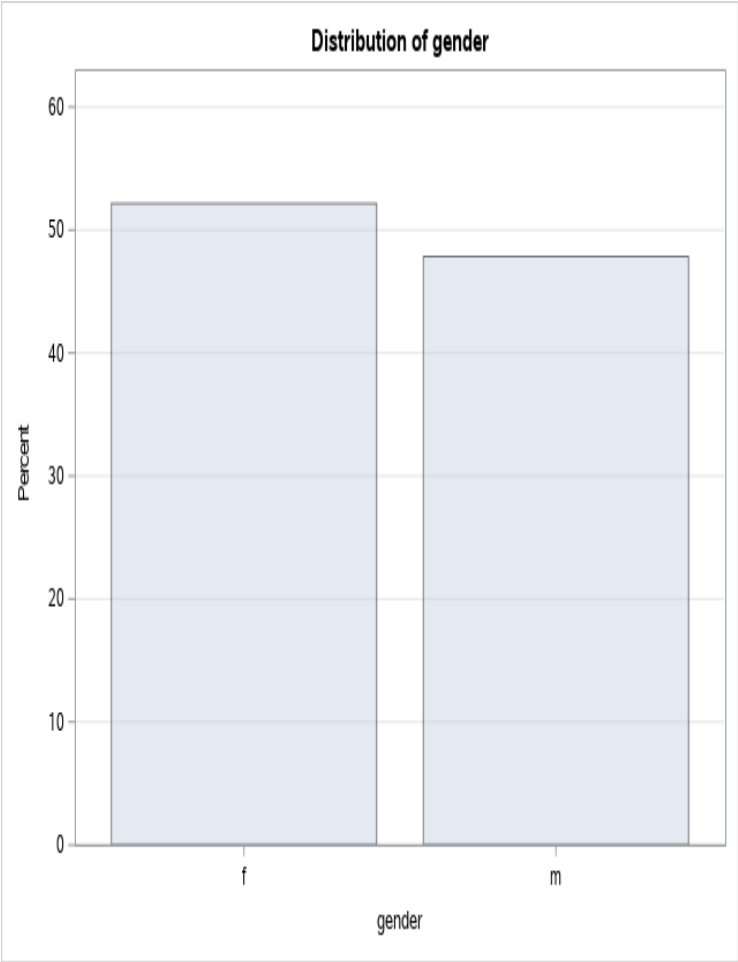
Relative frequency distribution for gender

```
In [8]: proc freq data=work.cdc;
         tables gender / plots=freqplot(scale=percent);
         run;
```

**Out[8]:** **The SAS System**

**The FREQ Procedure**

gender	Frequency	Percent	Cumulative Frequency	Cumulative Percent
f	10431	52.16	10431	52.16
m	9569	47.85	20000	100.00



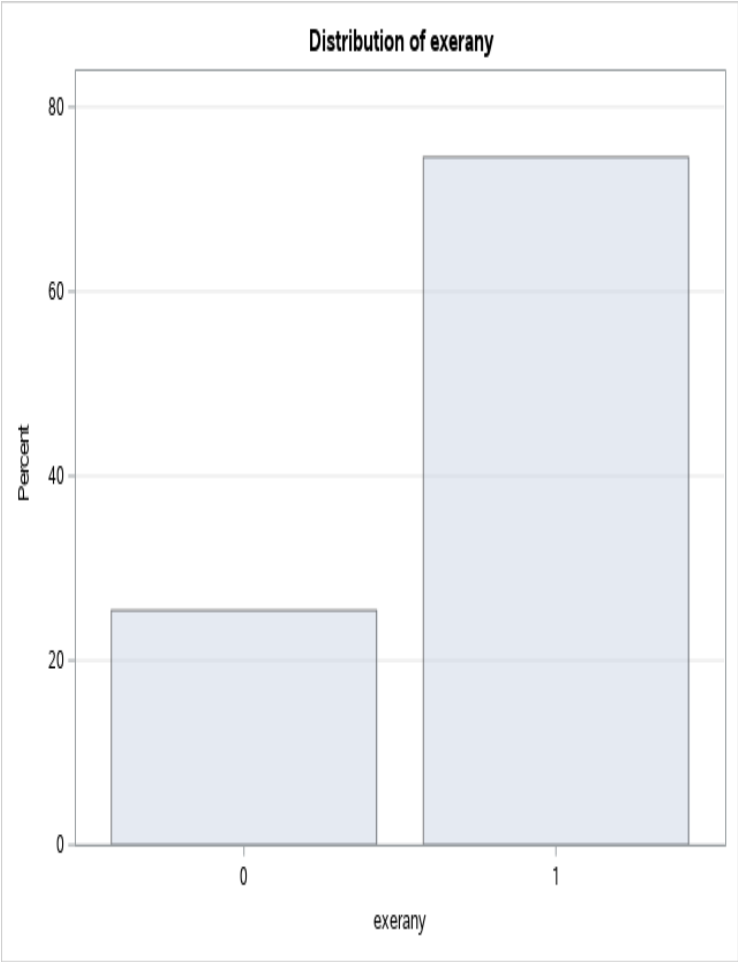
Relative frequency distribution for exerany

```
In [9]: proc freq data=work.cdc;
        tables exerany / plots=freqplot(scale=per
cent);
run;
```

Out[9]: **The SAS System**  
**The FREQ Procedure**

exerany	Frequency	Percent	Cumulative Frequency	Cumulative Percent
---------	-----------	---------	----------------------	--------------------

exerany	Frequency	Percent	Cumulative Frequency	C P
0	5086	25.43	5086	2
1	14914	74.57	20000	1



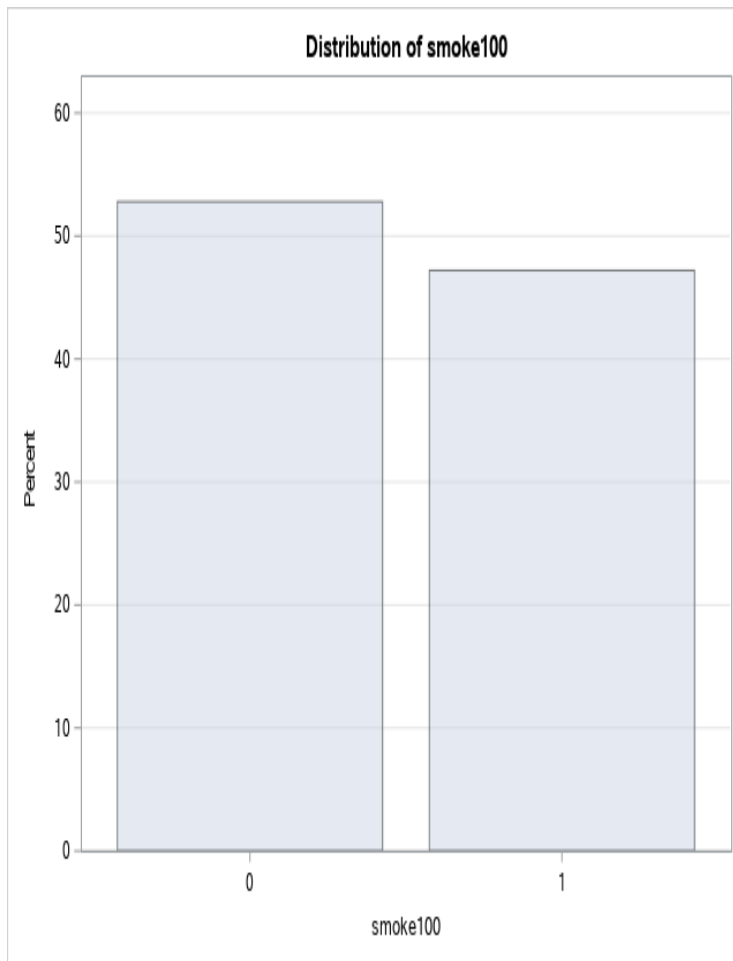
Relative frequency distribution for smoke100

```
In [10]: proc freq data=work.cdc;
          tables smoke100 / plots=freqplot(scale=pe
rcent);
run;
```

Out[10]: **The SAS System**  
**The FREQ Procedure**

smoke100	Frequency	Percent	Cumulative Frequency
----------	-----------	---------	-------------------------

smoke100	Frequency	Percent	Cumulative Frequency
0	10559	52.80	10559
1	9441	47.21	20000



Relative frequency distribution for genhlth

```
In [11]: proc freq data=work.cdc;
          tables genhlth / plots=freqplot(scale=per
cent);
run;
```

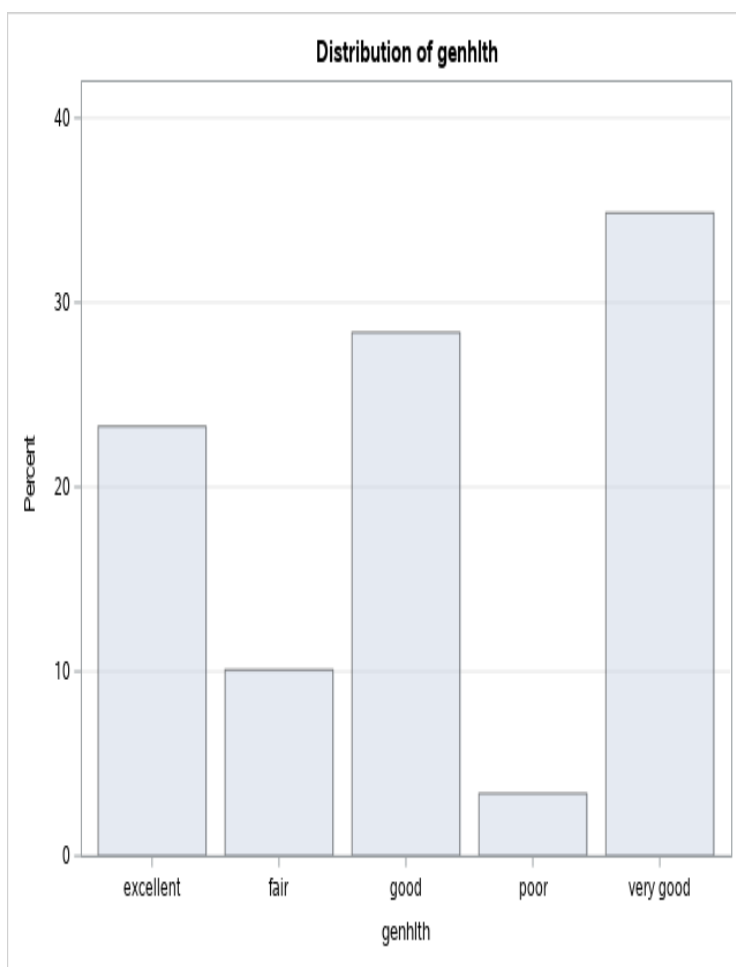
Out[11]: **The SAS System**

**The FREQ Procedure**

genhlth	Frequency	Percent	Cumulative Frequency
---------	-----------	---------	----------------------



genhlth	Frequency	Percent	Cumulative Frequency
excellent	4657	23.29	4657
fair	2019	10.10	6676
good	5675	28.38	12351
poor	677	3.39	13028
very good	6972	34.86	20000



Create multi-way frequency tables: gender and smoke100

```
In [12]: proc freq data=work.cdc;
          tables gender*smoke100;
          run;
```

Out[12]:

The SAS System

The FREQ Procedure

<div>Frequency Percent Row Pct Col Pct</div>	Table of gender by smoke100			
	gender	smoke100		
		0	1	Total
	f	6012	4419	10431
		30.06	22.10	52.16
		57.64	42.36	
		56.94	46.81	
	m	4547	5022	9569
		22.74	25.11	47.85
		47.52	52.48	
		43.06	53.19	
	Total	10559	9441	20000
		52.80	47.21	100.00

Create a mosaic plot

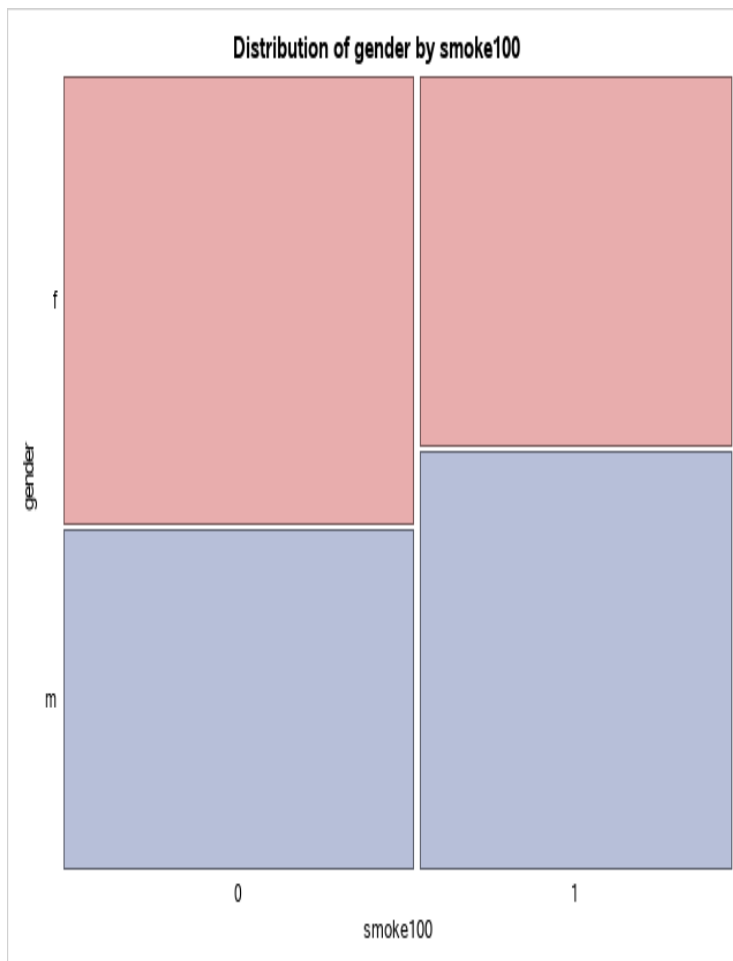
```
In [13]: proc freq data=work.cdc;
          tables Gender*Smoke100 / plots=mosaicplo
t;
run;
```

Out[13]: The SAS System

The FREQ Procedure

<div>Frequency Percent Row Pct Col Pct</div>	Table of gender by smoke100			
	gender	smoke100		
		0	1	Total
	f	6012	4419	10431
		30.06	22.10	52.16
		57.64	42.36	
		56.94	46.81	
	m	4547	5022	9569
		22.74	25.11	47.85
		47.52	52.48	
		43.06	53.19	
	Total	10559	9441	20000
		52.80	47.21	100.00

Table of gender by smoke100			
gender	smoke100		
	0	1	Total
m	4547	5022	9569
	22.74	25.11	47.85
	47.52	52.48	
	43.06	53.19	
Total	10559	9441	20000
	52.80	47.21	100.00



### Exercise 3

What does the mosaic plot reveal about smoking habits and gender?

The percentage of males that smoked at least 100 cigarettes in their entire life is larger than percentage of females

## Interlude: How SAS Processes Data

See the first 10 values in the data portion

```
In [14]: proc print data=work.cdc (obs=10);
run;
```

Out[14]:

The SAS System

Obs	genhlth	exerany	hlthplan	smoke100
1	good	0	1	0
2	good	0	1	1
3	good	1	1	1
4	good	1	1	0
5	very good	0	1	0
6	very good	1	1	0
7	very good	1	1	0
8	very good	0	1	0
9	good	0	1	1
10	good	1	1	0

See the descriptor portion such as the names, types, and lengths of the variables

```
In [15]: proc contents data=work.cdc;
run;
```

Out[15]:

The SAS System

## The CONTENTS Procedure

<b>Data Set Name</b>	WORK.CDC	<b>Observations</b>
<b>Member Type</b>	DATA	<b>Variables</b>
<b>Engine</b>	V9	<b>Indexes</b>
<b>Created</b>	07/07/2018 06:06:24	<b>Observations Length</b>
<b>Last Modified</b>	07/07/2018 06:06:24	<b>Deleted Observations</b>
<b>Protection</b>		<b>Compressed</b>
<b>Data Set Type</b>		<b>Sorted</b>
<b>Label</b>		
<b>Data Representation</b>	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64	
<b>Encoding</b>	utf-8 Unicode (UTF-8)	

### Engine/Host Dependent Information

<b>Data Set Page Size</b>	65536
<b>Number of Data Set Pages</b>	23
<b>First Data Page</b>	1
<b>Max Obs per Page</b>	908

### Engine/Host Dependent Information

<b>Obs in First Data Page</b>	867
<b>Number of Data Set Repairs</b>	0
<b>Filename</b>	/tmp/SAS_work89DD00004797_loca
<b>Release Created</b>	9.0401M5
<b>Host Created</b>	Linux
<b>Inode Number</b>	409728
<b>Access Permission</b>	rw-r--r--
<b>Owner Name</b>	sasdemo
<b>File Size</b>	2MB
<b>File Size (bytes)</b>	1572864

### Alphabetic List of Variables and Attributes

#	Variable	Type	Len	Format	Informat
8	age	Num	8	BEST12.	BEST32.
2	exerany	Num	8	BEST12.	BEST32.
9	gender	Char	1	\$1.	\$1.
1	genhlth	Char	9	\$9.	\$9.
5	height	Num	8	BEST12.	BEST32.

### Alphabetic List of Variables and Attributes

#	Variable	Type	Len	Format	Informat
3	hlthplan	Num	8	BEST12.	BEST32.
4	smoke100	Num	8	BEST12.	BEST32.
6	weight	Num	8	BEST12.	BEST32.
7	wtdesire	Num	8	BEST12.	BEST32.

Create subset of observations of people who are men or anyone over the age of 30. See the first 10 values of work.newcdc

```
In [16]: data work.newcdc;
          set work.cdc;
          if gender="m" and age>30;
run;

proc print data=work.newcdc (obs=10);
run;
```

Out[16]: **The SAS System**

Obs	genhlth	exerany	hlthplan	smoke100
1	good	0	1	0
2	very good	1	1	0
3	very good	0	1	0
4	good	1	1	0
5	excellent	1	1	1
6	fair	1	1	1
7	excellent	1	1	1
8	good	1	1	1

Obs	genhlth	exerany	hlthplan	smoke100
9	good	0	0	1
10	fair	0	1	1

Create subset of observations of people who are men or female and less the age of 30. See the first 10 values of work.newcdc

```
In [17]: data work.newcdc;
          set work.cdc;
          if gender="m" or (gender="f" and age>30);
run;

proc print data=work.newcdc (obs=10);
run;
```

Out[17]:

The SAS System

Obs	genhlth	exerany	hlthplan	smoke100
1	good	0	1	0
2	good	0	1	1
3	good	1	1	1
4	good	1	1	0
5	very good	0	1	0
6	very good	1	1	0
7	very good	1	1	0
8	very good	0	1	0
9	good	1	1	0
10	excellent	1	1	1



## Exercise 4

Create a new data set named under23smoke that contains all observations of respondents under the age of 23 that have smoked 100 cigarettes in their lifetime. Write the programming statements you used to create the new data set as the answer to this exercise.

```
In [18]: data work.under23smoke;
          set work.cdc;
          if smoke100=1 and age<23;
        run;

        proc print data=work.under23smoke (obs=10);
        run;
```

Out[18]:

The SAS System

Obs	genhlth	exerany	hlthplan	smoke100
1	excellent	1	0	1
2	very good	1	0	1
3	excellent	1	1	1
4	good	1	1	1
5	very good	1	1	1
6	very good	1	0	1
7	fair	0	1	1
8	fair	1	1	1
9	excellent	1	0	1
10	fair	1	1	1

## Quantitative Data

Create box-and-whisker plot and a histogram for weight

```
In [19]: ods graphics;
proc univariate data=work.cdc plots;
    var weight;
run;
```

Out[19]:

### The SAS System

#### The UNIVARIATE Procedure Variable: weight

Moments			
<b>N</b>	20000	<b>Sum Weights</b>	20000
<b>Mean</b>	169.68295	<b>Sum Observations</b>	3393659
<b>Std Deviation</b>	40.08097	<b>Variance</b>	1606
<b>Skewness</b>	0.95572799	<b>Kurtosis</b>	1.99
<b>Uncorrected SS</b>	607974147	<b>Corrected SS</b>	3212
<b>Coeff Variation</b>	23.6210945	<b>Std Error Mean</b>	0.28

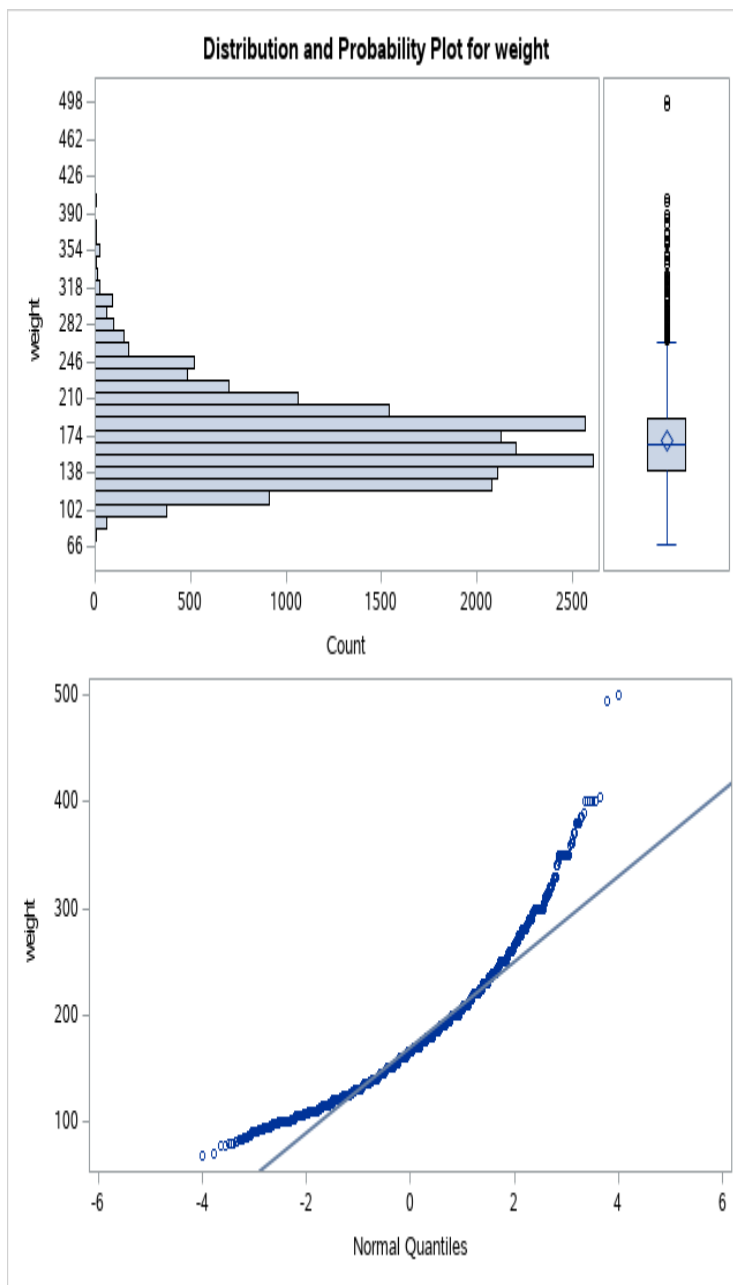
Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	169.6830	<b>Std Deviation</b>	40.08097
<b>Median</b>	165.0000	<b>Variance</b>	1606
<b>Mode</b>	160.0000	<b>Range</b>	432.00000
		<b>Interquartile Range</b>	50.00000

Tests for Location: Mu0=0				
Test	Statistic		p Value	
Student's t	t	598.7079	Pr >  t	<.0001
Sign	M	10000	Pr >=  M	<.0001
Signed Rank	S	1.0001E8	Pr >=  S	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	500
99%	290
95%	240
90%	220
75% Q3	190
50% Median	165
25% Q1	140
10%	124
5%	115
1%	100
0% Min	68

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
68	18743	400	2944
70	16531	400	19319
78	18065	405	15720
78	11299	495	4445
79	7614	500	1995



Create box-and-whisker plot and a histogram for weight subset by gender using "class" statement

```
In [20]: ods graphics;
proc univariate data=work.cdc plots;
  class gender;
  var weight;
run;
```

Out[20]: **The SAS System**

**The UNIVARIATE Procedure**  
**Variable: weight**  
**gender = f**

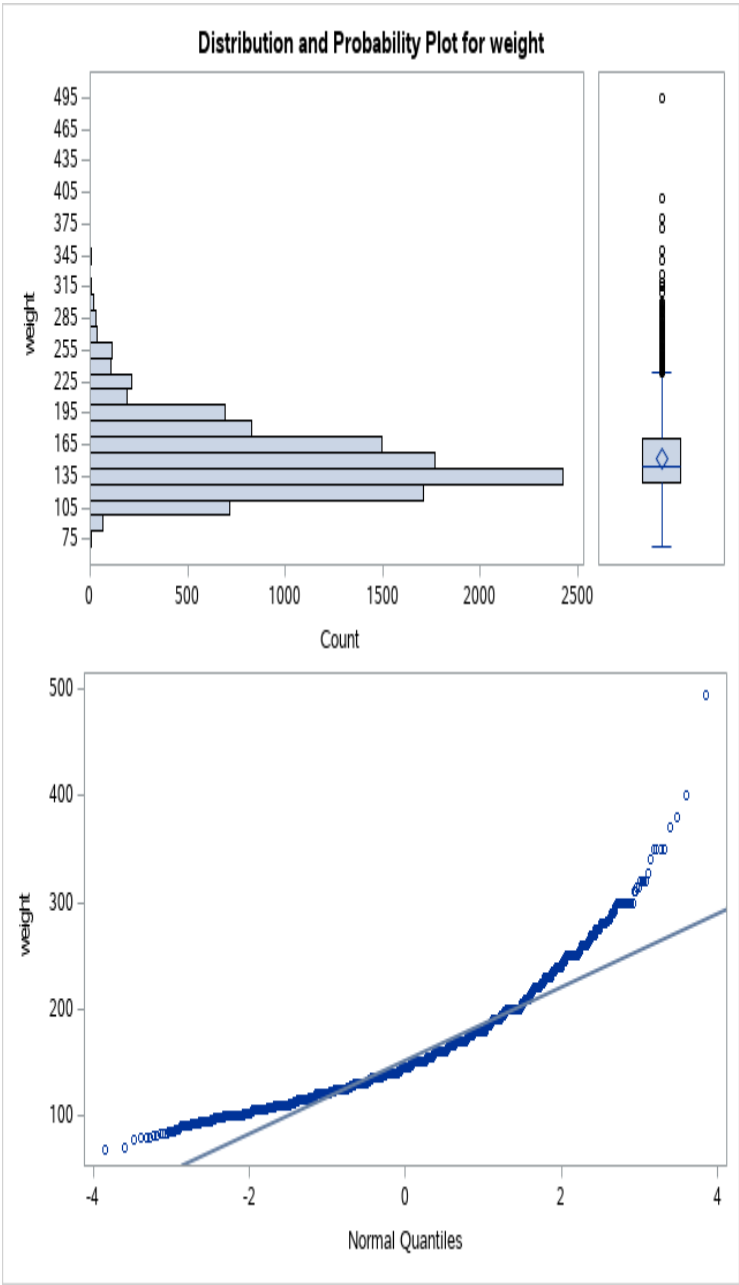
<b>Moments</b>			
<b>N</b>	10431	<b>Sum Weights</b>	10431
<b>Mean</b>	151.666187	<b>Sum Observations</b>	1582000
<b>Std Deviation</b>	34.2975191	<b>Variance</b>	1176.00000
<b>Skewness</b>	1.38530136	<b>Kurtosis</b>	3.77
<b>Uncorrected SS</b>	252209474	<b>Corrected SS</b>	122600000
<b>Coeff Variation</b>	22.6138203	<b>Std Error Mean</b>	0.33

<b>Basic Statistical Measures</b>			
<b>Location</b>		<b>Variability</b>	
<b>Mean</b>	151.6662	<b>Std Deviation</b>	34.29752
<b>Median</b>	145.0000	<b>Variance</b>	1176
<b>Mode</b>	140.0000	<b>Range</b>	427.00000
		<b>Interquartile Range</b>	42.00000

<b>Tests for Location: <math>\mu_0=0</math></b>				
<b>Test</b>	<b>Statistic</b>		<b>p Value</b>	
<b>Student's t</b>	<b>t</b>	451.6365	<b>Pr &gt;  t </b>	<.0001
<b>Sign</b>	<b>M</b>	5215.5	<b>Pr &gt;=  M </b>	<.0001
<b>Signed Rank</b>	<b>S</b>	27204048	<b>Pr &gt;=  S </b>	<.0001

<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>
<b>100% Max</b>	495
<b>99%</b>	260
<b>95%</b>	220
<b>90%</b>	198
<b>75% Q3</b>	170
<b>50% Median</b>	145
<b>25% Q1</b>	128
<b>10%</b>	115
<b>5%</b>	110
<b>1%</b>	100
<b>0% Min</b>	68

<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
68	18743	350	13607
70	16531	371	4612
78	11299	380	7160
79	7614	400	19319
80	15673	495	4445



**The SAS System**

**The UNIVARIATE Procedure**  
**Variable: weight**  
**gender = m**

Moments			
N	9569	Sum Weights	9569
Mean	189.322709	Sum Observations	1811



<b>Moments</b>			
<b>Std Deviation</b>	36.5503551	<b>Variance</b>	1336
<b>Skewness</b>	1.16960127	<b>Kurtosis</b>	3.04
<b>Uncorrected SS</b>	355764673	<b>Corrected SS</b>	1278
<b>Coeff Variation</b>	19.3058484	<b>Std Error Mean</b>	0.37

<b>Basic Statistical Measures</b>			
<b>Location</b>		<b>Variability</b>	
<b>Mean</b>	189.3227	<b>Std Deviation</b>	36.55036
<b>Median</b>	185.0000	<b>Variance</b>	1336
<b>Mode</b>	180.0000	<b>Range</b>	422.00000
		<b>Interquartile Range</b>	45.00000

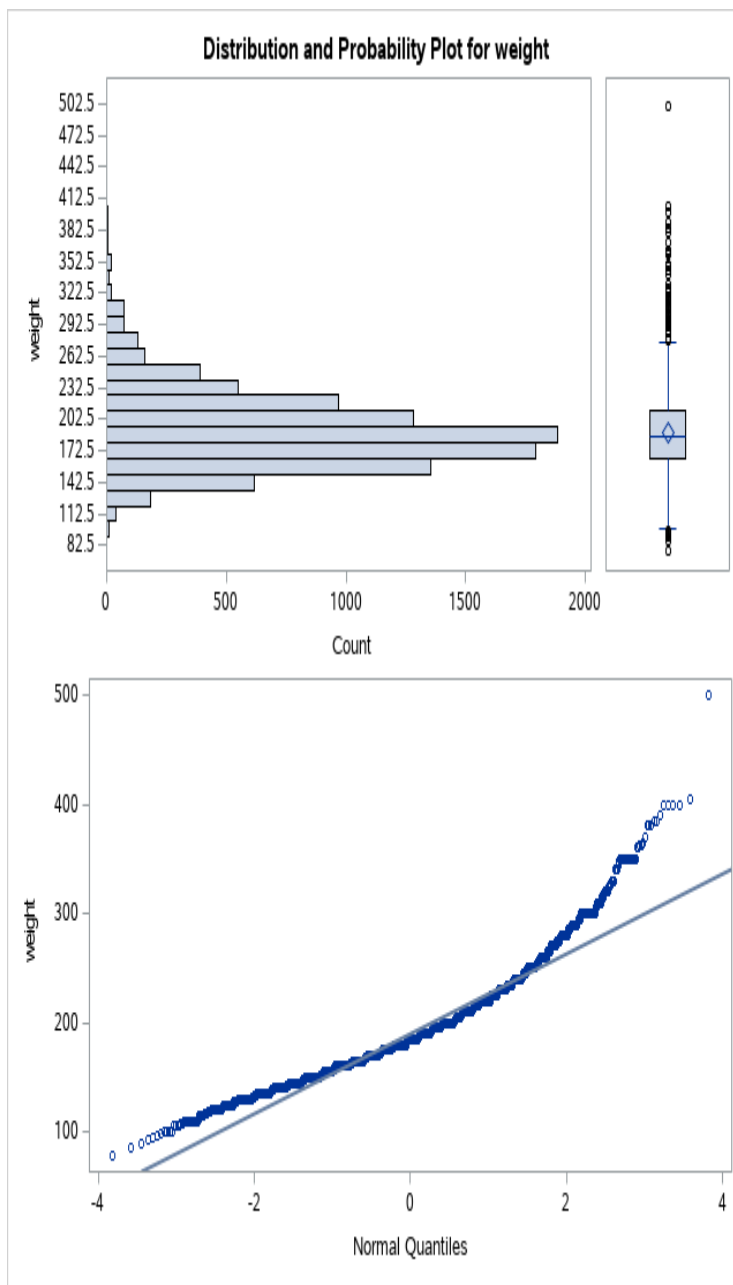
<b>Tests for Location: <math>\mu_0=0</math></b>				
<b>Test</b>	<b>Statistic</b>		<b>p Value</b>	
<b>Student's t</b>	<b>t</b>	506.6924	<b>Pr &gt;  t </b>	<.0001
<b>Sign</b>	<b>M</b>	4784.5	<b>Pr &gt;=  M </b>	<.0001
<b>Signed Rank</b>	<b>S</b>	22893833	<b>Pr &gt;=  S </b>	<.0001

<b>Quantiles (Definition 5)</b>
---------------------------------

<b>Quantiles (Definition 5)</b>	<b>Quantile</b>
---------------------------------	-----------------

<b>Level</b>	<b>Quantile</b>
<b>100% Max</b>	500
<b>99%</b>	300
<b>95%</b>	256
<b>90%</b>	235
<b>75% Q3</b>	210
<b>50% Median</b>	185
<b>25% Q1</b>	165
<b>10%</b>	150
<b>5%</b>	140
<b>1%</b>	125
<b>0% Min</b>	78

<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
78	18065	400	1279
86	15967	400	2659
90	303	400	2944
93	9558	405	15720
94	11333	500	1995



Create box-and-whisker plot and a histogram for weight subset by gender using "by" statement

```
In [21]: proc sort data=work.cdc;
          by gender;
          run;

          ods graphics;
          proc univariate data=work.cdc plots;
            by gender;
            var Weight;
          run;
```

Out[21]: **The SAS System**

**The UNIVARIATE Procedure**

Variable: weight

gender=f

Moments			
<b>N</b>	10431	<b>Sum Weights</b>	10431
<b>Mean</b>	151.666187	<b>Sum Observations</b>	1582094.74
<b>Std Deviation</b>	34.2975191	<b>Variance</b>	1176.297519
<b>Skewness</b>	1.38530136	<b>Kurtosis</b>	3.770136
<b>Uncorrected SS</b>	252209474	<b>Corrected SS</b>	122629.74
<b>Coeff Variation</b>	22.6138203	<b>Std Error Mean</b>	0.333333

Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	151.6662	<b>Std Deviation</b>	34.29752
<b>Median</b>	145.0000	<b>Variance</b>	1176
<b>Mode</b>	140.0000	<b>Range</b>	427.00000
		<b>Interquartile Range</b>	42.00000

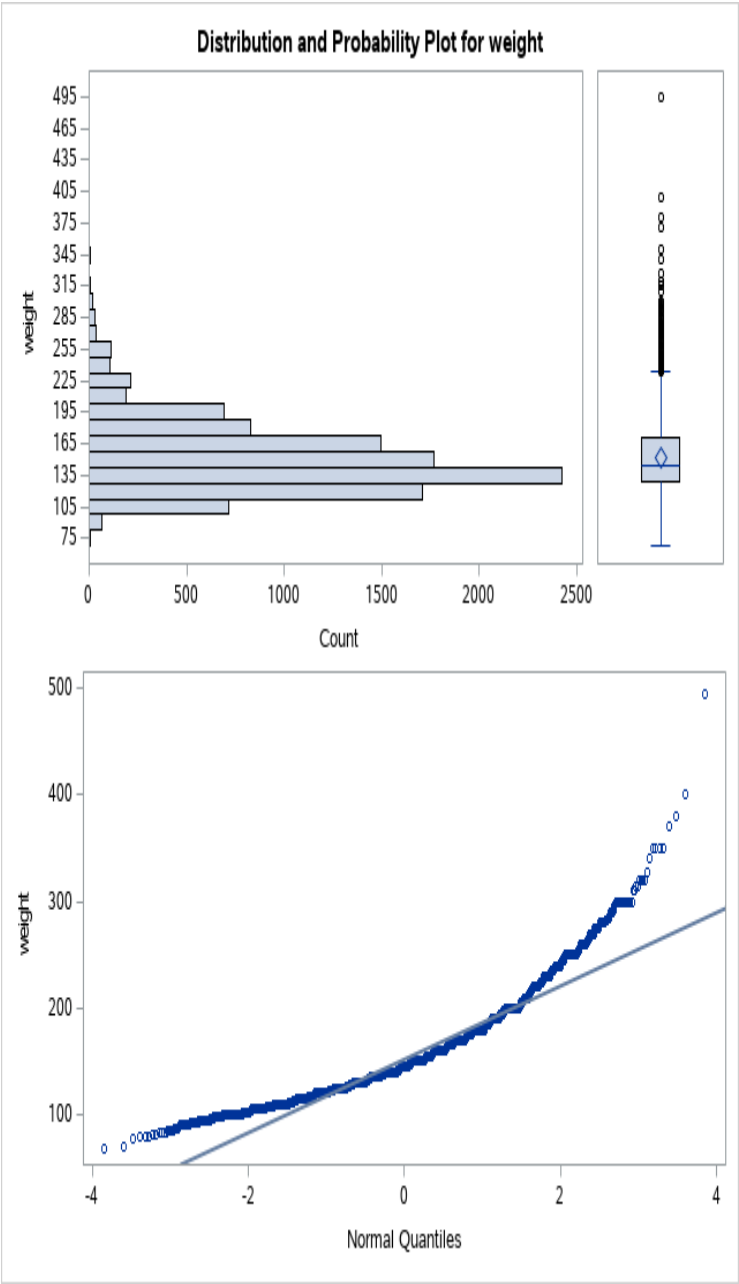
Tests for Location: Mu0=0				
Test	Statistic		p Value	
<b>Student's t</b>	<b>t</b>	451.6365	<b>Pr &gt;  t </b>	<.0001

Tests for Location: $\mu_0=0$				
Test	Statistic		p Value	
Sign	M	5215.5	Pr $\geq  M $	<.0001
Signed Rank	S	27204048	Pr $\geq  S $	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	495
99%	260
95%	220
90%	198
75% Q3	170
50% Median	145
25% Q1	128
10%	115
5%	110
1%	100
0% Min	68

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
68	9739	350	7051

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
70	8592	371	2354
78	5860	380	3663
79	3905	400	10060
80	8154	495	2284



## The SAS System

### The UNIVARIATE Procedure

Variable: weight

gender=m

Moments			
<b>N</b>	9569	<b>Sum Weights</b>	9569
<b>Mean</b>	189.322709	<b>Sum Observations</b>	1811
<b>Std Deviation</b>	36.5503551	<b>Variance</b>	1335
<b>Skewness</b>	1.16960127	<b>Kurtosis</b>	3.04
<b>Uncorrected SS</b>	355764673	<b>Corrected SS</b>	1278
<b>Coeff Variation</b>	19.3058484	<b>Std Error Mean</b>	0.37

Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	189.3227	<b>Std Deviation</b>	36.55036
<b>Median</b>	185.0000	<b>Variance</b>	1336
<b>Mode</b>	180.0000	<b>Range</b>	422.00000
		<b>Interquartile Range</b>	45.00000

Tests for Location: Mu0=0		
Test	Statistic	p Value

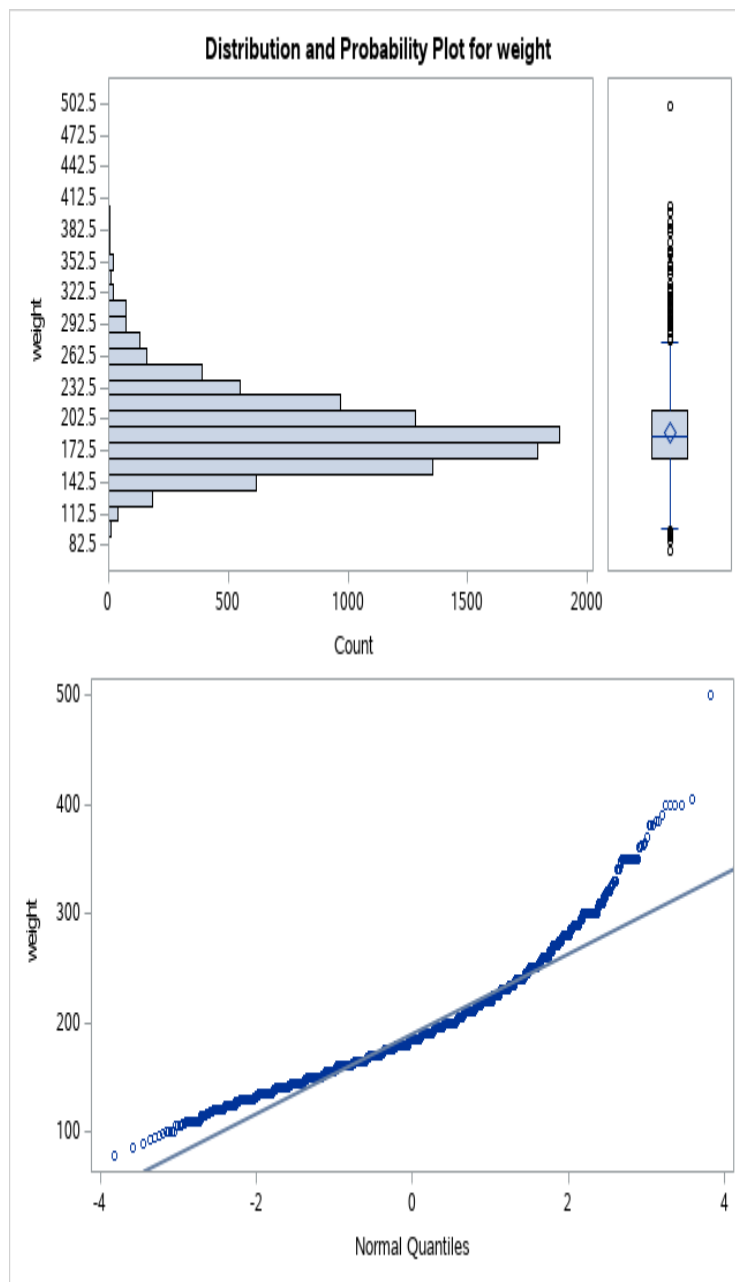
Tests for Location: Mu0=0				
Test	Statistic		p Value	
Student's t	t	506.6924	Pr >  t	<.0001
Sign	M	4784.5	Pr >=  M	<.0001
Signed Rank	S	22893833	Pr >=  S	<.0001

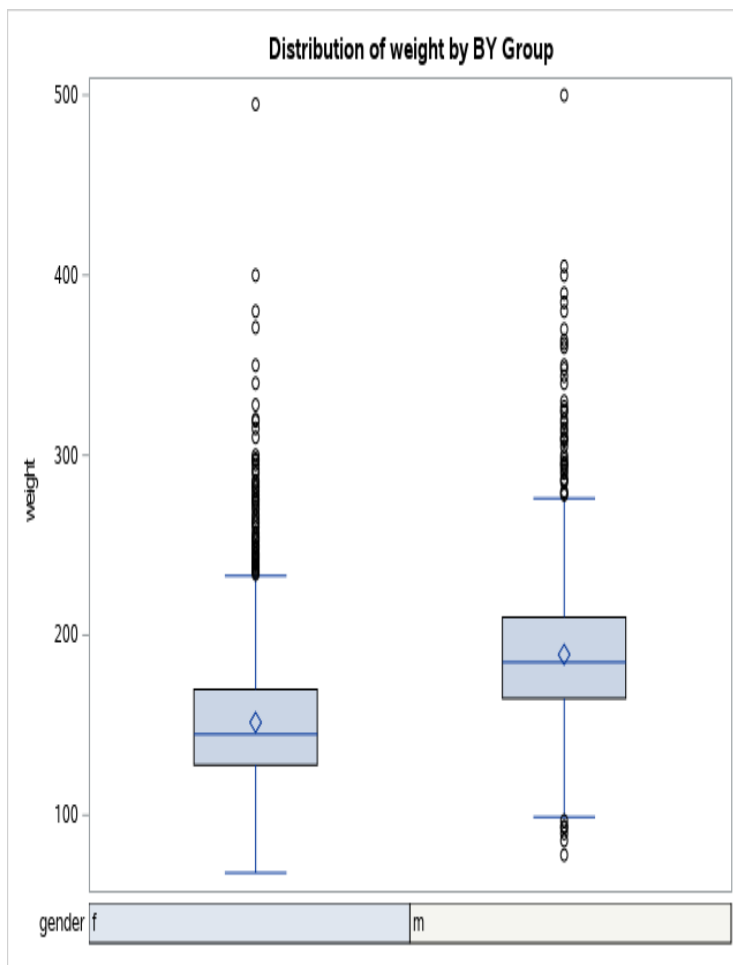
Quantiles (Definition 5)	
Level	Quantile
100% Max	500
99%	300
95%	256
90%	235
75% Q3	210
50% Median	185
25% Q1	165
10%	150
5%	140
1%	125
0% Min	78

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs



Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
78	19127	400	11036
86	18099	400	11725
90	10578	400	11867
93	15074	405	17968
94	15887	500	11392





Data set with a new variable named bmi

```
In [22]: data work.cdcbmi;
         set work.cdc;
         bmi = (weight / height**2) * 703;
         run;
```

Out[22]:

```
288 ods listing close;ods html5 (id=saspy_
internal) file=stdout options(bitmap_mode
='inline') device=svg; ods graphics on /
288! outputfmt=png;
NOTE: Writing HTML5(SASPY_INTERNAL) Body fi
le: STDOUT
289
290 data work.cdcbmi;
291     set work.cdc;
292     bmi = (weight / height**2) * 703;
293 run;
NOTE: There were 20000 observations read fr
om the data set WORK.CDC.
NOTE: The data set WORK.CDCBMI has 20000 ob
servations and 10 variables.
NOTE: DATA statement used (Total process ti
me):
      real time          0.00 seconds
```

cpu time

0.00 seconds

```

294
295 ods html5 (id=saspy_internal) close;od
s listing;

296

```

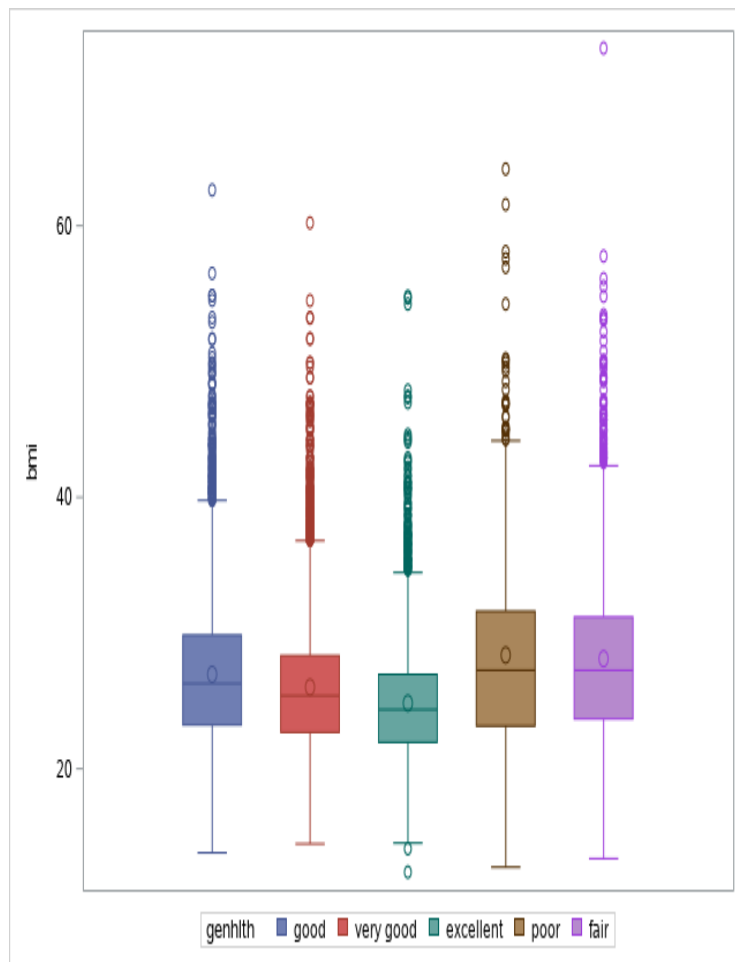
Alternative method to producing box-and-whisker plots using group statement

```

In [23]: proc sgplot data=work.cdcbmi;
          vbox bmi / group=genhlth;
          run;

```

Out[23]:



## Exercise 5

What does this box plot show? Pick another categorical variable from the data set and see how it relates to BMI. List the variable you chose, why you might think it would have a relationship to BMI, and indicate what the figure seems to suggest.

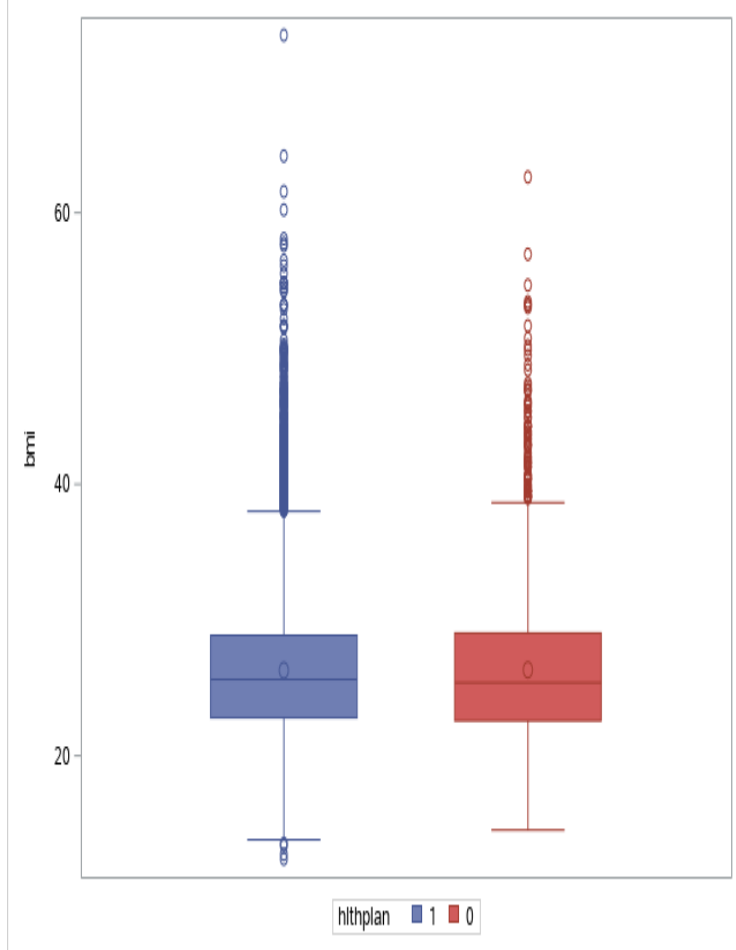
```

In [24]: proc sgplot data=work.cdcbmi;

```

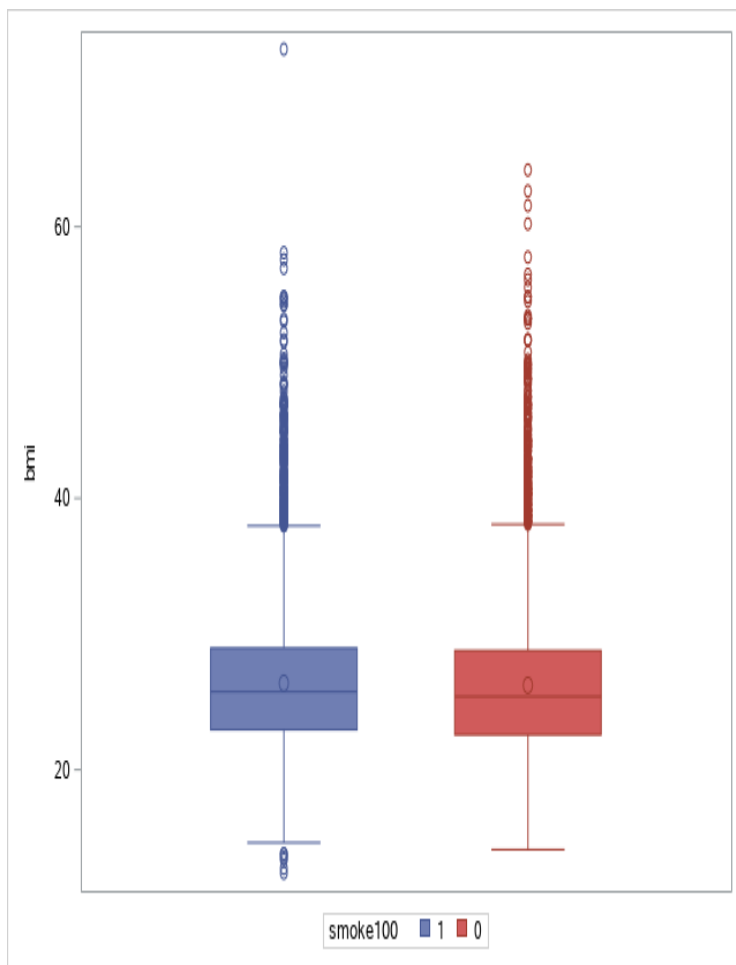
```
vbox bmi / group=hlthplan;  
run;
```

Out[24]:



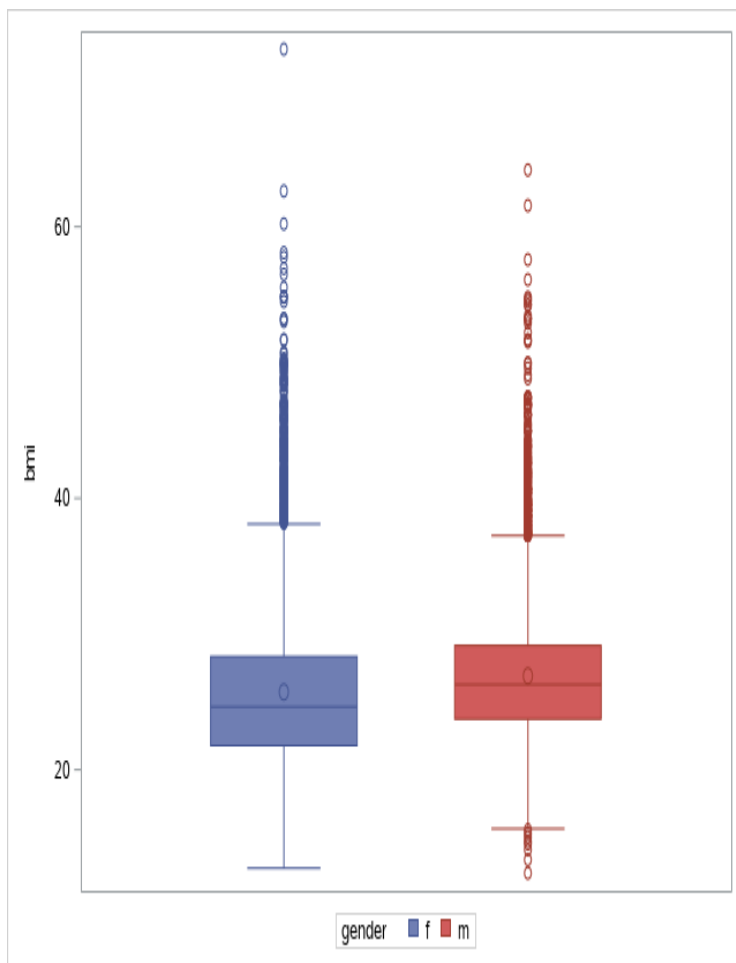
```
In [25]: proc sgplot data=work.cdcbmi;  
          vbox bmi / group=smoke100;  
run;
```

Out[25]:



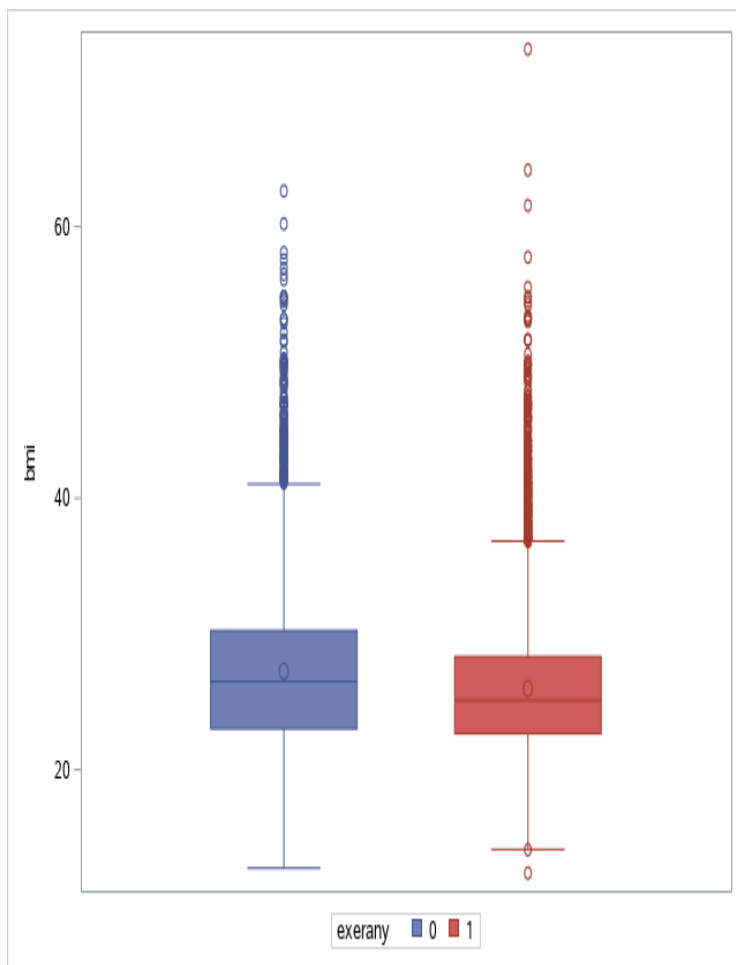
```
In [26]: proc sgplot data=work.cdcbmi;  
          vbox bmi / group=gender;  
          run;
```

Out[26]:



```
In [27]: proc sgplot data=work.cdcbmi;  
          vbox bmi / group=exerany;  
          run;
```

Out[27]:



Histogram for the bmi of the respondents

```
In [28]: ods graphics;
proc univariate data=work.cdcbmi;
  var bmi;
  histogram bmi;
run;
```

Out[28]:

## The SAS System

The UNIVARIATE Procedure  
Variable: bmi

Moments			
N	20000	Sum Weights	20000
Mean	26.3069252	Sum Observations	526145.2
Std Deviation	5.21810488	Variance	27.2367



<b>Moments</b>			
<b>Skewness</b>	1.27589014	<b>Kurtosis</b>	3.32
<b>Uncorrected SS</b>	14385631.4	<b>Corrected SS</b>	5445
<b>Coeff Variation</b>	19.83548	<b>Std Error Mean</b>	0.03

<b>Basic Statistical Measures</b>			
<b>Location</b>		<b>Variability</b>	
<b>Mean</b>	26.30693	<b>Std Deviation</b>	5.21810
<b>Median</b>	25.60354	<b>Variance</b>	27.22862
<b>Mode</b>	27.12191	<b>Range</b>	60.69029
		<b>Interquartile Range</b>	6.17852

<b>Tests for Location: Mu0=0</b>				
<b>Test</b>	<b>Statistic</b>		<b>p Value</b>	
<b>Student's t</b>	<b>t</b>	712.9717	<b>Pr &gt;  t </b>	<.0001
<b>Sign</b>	<b>M</b>	10000	<b>Pr &gt;=  M </b>	<.0001
<b>Signed Rank</b>	<b>S</b>	1.0001E8	<b>Pr &gt;=  S </b>	<.0001

<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>
<b>100% Max</b>	73.0907

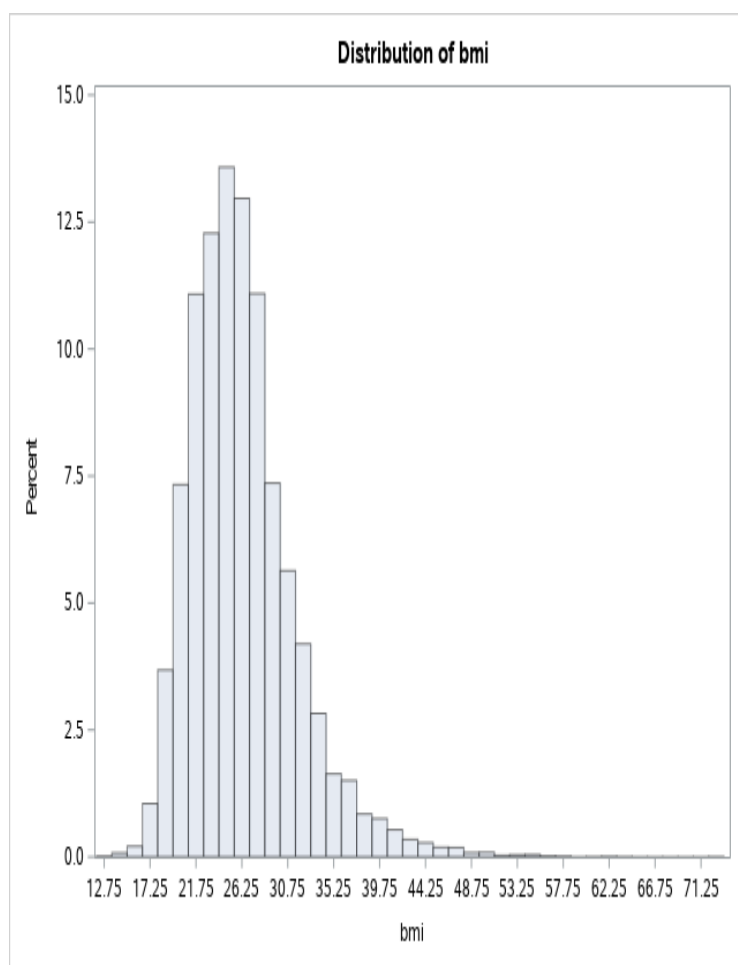
<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>
<b>99%</b>	43.5553
<b>95%</b>	35.9509
<b>90%</b>	32.8802
<b>75% Q3</b>	28.8862
<b>50% Median</b>	25.6035
<b>25% Q1</b>	22.7077
<b>10%</b>	20.5957
<b>5%</b>	19.4835
<b>1%</b>	17.7123
<b>0% Min</b>	12.4005

<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
12.4005	15887	60.2217	443
12.7495	3905	61.5733	17968
13.3872	19127	62.6420	10060
13.5565	8698	64.1892	11392
13.8104	2493	73.0907	2284

---

**The SAS System**

**The UNIVARIATE Procedure**



Histogram for the bmi of the respondents grouped by gender

```
In [29]: ods graphics;
proc univariate data=work.cdcbmi;
  class gender;
  var bmi;
  histogram bmi;
run;
```

Out[29]:

**The SAS System**

**The UNIVARIATE Procedure**

**Variable: bmi**

**gender = f**

Moments			
<b>N</b>	10431	<b>Sum Weights</b>	10431
<b>Mean</b>	25.7411474	<b>Sum Observations</b>	268500

<b>Moments</b>			
<b>Std Deviation</b>	5.62057729	<b>Variance</b>	31.5
<b>Skewness</b>	1.36746824	<b>Kurtosis</b>	3.13
<b>Uncorrected SS</b>	7241143.13	<b>Corrected SS</b>	3294
<b>Coeff Variation</b>	21.8349913	<b>Std Error Mean</b>	0.05

<b>Basic Statistical Measures</b>			
<b>Location</b>		<b>Variability</b>	
<b>Mean</b>	25.74115	<b>Std Deviation</b>	5.62058
<b>Median</b>	24.63832	<b>Variance</b>	31.59089
<b>Mode</b>	27.46094	<b>Range</b>	60.34120
		<b>Interquartile Range</b>	6.55244

<b>Tests for Location: Mu0=0</b>				
<b>Test</b>	<b>Statistic</b>		<b>p Value</b>	
<b>Student's t</b>	<b>t</b>	467.7459	<b>Pr &gt;  t </b>	<.0001
<b>Sign</b>	<b>M</b>	5215.5	<b>Pr &gt;=  M </b>	<.0001
<b>Signed Rank</b>	<b>S</b>	27204048	<b>Pr &gt;=  S </b>	<.0001

<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>

<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>
<b>100% Max</b>	73.0907
<b>99%</b>	44.3980
<b>95%</b>	36.5765
<b>90%</b>	32.9453
<b>75% Q3</b>	28.3396
<b>50% Median</b>	24.6383
<b>25% Q1</b>	21.7872
<b>10%</b>	19.8529
<b>5%</b>	18.9853
<b>1%</b>	17.2684
<b>0% Min</b>	12.7495

<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
12.7495	3905	57.7725	3663
13.5565	8698	58.1005	2354
13.8104	2493	60.2217	443
13.8156	5860	62.6420	10060
14.1367	8592	73.0907	2284

---

**The SAS System**

**The UNIVARIATE Procedure**

Variable: bmi  
gender = m

Moments			
<b>N</b>	9569	<b>Sum Weights</b>	9569
<b>Mean</b>	26.9236698	<b>Sum Observations</b>	257600
<b>Std Deviation</b>	4.6633462	<b>Variance</b>	21.7357
<b>Skewness</b>	1.30603659	<b>Kurtosis</b>	3.92
<b>Uncorrected SS</b>	7144488.31	<b>Corrected SS</b>	208000
<b>Coeff Variation</b>	17.3206188	<b>Std Error Mean</b>	0.0477

Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	26.92367	<b>Std Deviation</b>	4.66335
<b>Median</b>	26.28313	<b>Variance</b>	21.74680
<b>Mode</b>	27.12191	<b>Range</b>	51.78874
		<b>Interquartile Range</b>	5.41028

Tests for Location: Mu0=0				
Test	Statistic		p Value	
<b>Student's t</b>	<b>t</b>	564.7677	<b>Pr &gt;  t </b>	<.0001
<b>Sign</b>	<b>M</b>	4784.5	<b>Pr &gt;=  M </b>	<.0001

Tests for Location: $\mu_0=0$				
Test	Statistic		p Value	
Signed Rank	S	22893833	Pr $\geq  S $	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	64.1892
99%	42.0390
95%	35.4379
90%	32.5463
75% Q3	29.1561
50% Median	26.2831
25% Q1	23.7458
10%	21.9247
5%	20.8031
1%	18.6510
0% Min	12.4005

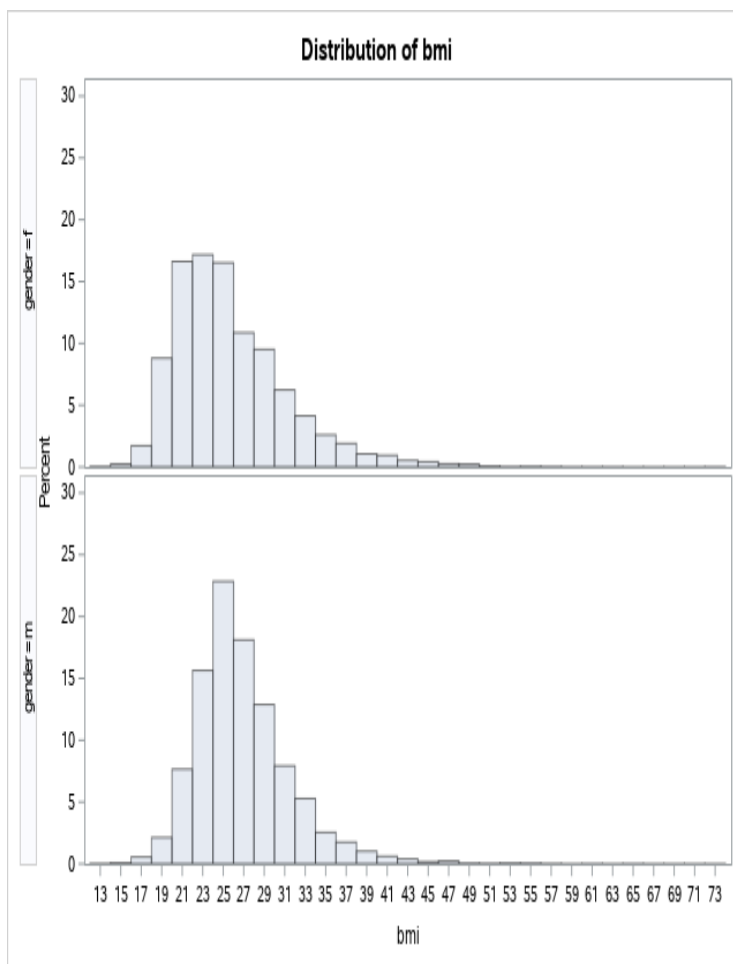
Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
12.4005	15887	54.8118	17510
13.3872	19127	56.1101	10796
14.1216	11534	57.5866	16778

## Extreme Observations

Lowest		Highest	
Value	Obs	Value	Obs
14.5493	18866	61.5733	17968
14.7603	18099	64.1892	11392

## The SAS System

### The UNIVARIATE Procedure



Histogram for the age of the respondents with 50 bins

```
In [30]: proc univariate data=work.cdcbmi;
          var bmi;
          histogram bmi / nmidpoints=50;
          run;
```



Out[30]:

## The SAS System

## The UNIVARIATE Procedure

Variable: bmi

Moments			
N	20000	Sum Weights	20000
Mean	26.3069252	Sum Observations	526138.56
Std Deviation	5.21810488	Variance	27.2378
Skewness	1.27589014	Kurtosis	3.3211
Uncorrected SS	14385631.4	Corrected SS	5445.12
Coeff Variation	19.83548	Std Error Mean	0.0352

Basic Statistical Measures			
Location		Variability	
Mean	26.30693	Std Deviation	5.21810
Median	25.60354	Variance	27.22862
Mode	27.12191	Range	60.69029
		Interquartile Range	6.17852

Tests for Location: Mu0=0				
Test	Statistic		p Value	
Student's t	t	712.9717	Pr >  t	<.0001

Tests for Location: $\mu_0=0$				
Test	Statistic		p Value	
Sign	M	10000	Pr $\geq$  M	<.0001
Signed Rank	S	1.0001E8	Pr $\geq$  S	<.0001

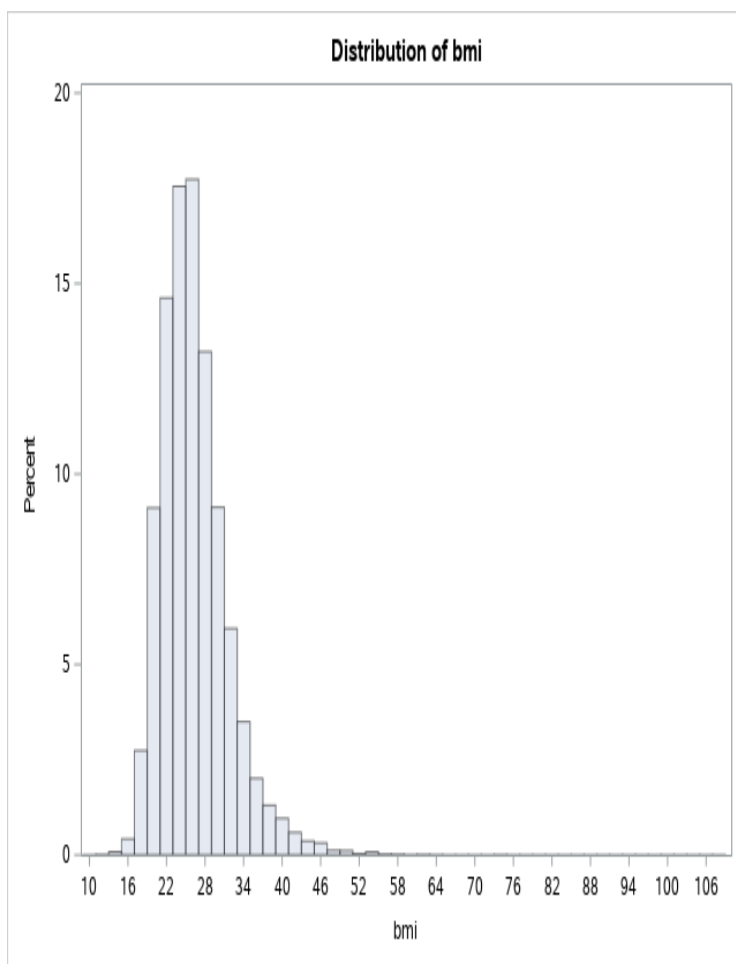
Quantiles (Definition 5)	
Level	Quantile
100% Max	73.0907
99%	43.5553
95%	35.9509
90%	32.8802
75% Q3	28.8862
50% Median	25.6035
25% Q1	22.7077
10%	20.5957
5%	19.4835
1%	17.7123
0% Min	12.4005

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
12.4005	15887	60.2217	443

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
12.7495	3905	61.5733	17968
13.3872	19127	62.6420	10060
13.5565	8698	64.1892	11392
13.8104	2493	73.0907	2284

### The SAS System

#### The UNIVARIATE Procedure



Histogram for the age of the respondents with 50 bins include the mean, standard deviation, and median

```
In [31]: ods graphics;
proc univariate data=work.cdcbmi;
  var bmi;
  histogram bmi / nmidpoints=50;
  inset mean std median / position=NE;
run;
```

Out[31]:

**The SAS System****The UNIVARIATE Procedure**

Variable: bmi

<b>Moments</b>			
<b>N</b>	20000	<b>Sum Weights</b>	20000
<b>Mean</b>	26.3069252	<b>Sum Observations</b>	526139.6
<b>Std Deviation</b>	5.21810488	<b>Variance</b>	27.2381
<b>Skewness</b>	1.27589014	<b>Kurtosis</b>	3.3211
<b>Uncorrected SS</b>	14385631.4	<b>Corrected SS</b>	5445.0
<b>Coeff Variation</b>	19.83548	<b>Std Error Mean</b>	0.0343

<b>Basic Statistical Measures</b>			
<b>Location</b>		<b>Variability</b>	
<b>Mean</b>	26.30693	<b>Std Deviation</b>	5.21810
<b>Median</b>	25.60354	<b>Variance</b>	27.22862
<b>Mode</b>	27.12191	<b>Range</b>	60.69029
		<b>Interquartile Range</b>	6.17852

<b>Tests for Location: Mu0=0</b>
----------------------------------

Tests for Location	Statistic	Test Value	p Value
--------------------	-----------	------------	---------

Test	Statistic		p Value	
Student's t	t	712.9717	Pr >  t	<.0001
Sign	M	10000	Pr >=  M	<.0001
Signed Rank	S	1.0001E8	Pr >=  S	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	73.0907
99%	43.5553
95%	35.9509
90%	32.8802
75% Q3	28.8862
50% Median	25.6035
25% Q1	22.7077
10%	20.5957
5%	19.4835
1%	17.7123
0% Min	12.4005

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
12.4005	15887	60.2217	443
12.7495	3905	61.5733	17968
13.3872	19127	62.6420	10060
13.5565	8698	64.1892	11392
13.8104	2493	73.0907	2284

### The SAS System

#### The UNIVARIATE Procedure

