

Chapter 8: Reading Raw Data Files

8.1 Introduction to Reading Raw Data Files

8.2 Reading Standard Delimited Data

8.3 Reading Nonstandard Delimited Data

8.4 Handling Missing Data

Chapter 8: Reading Raw Data Files

8.1 Introduction to Reading Raw Data Files

8.2 Reading Standard Delimited Data

8.3 Reading Nonstandard Delimited Data

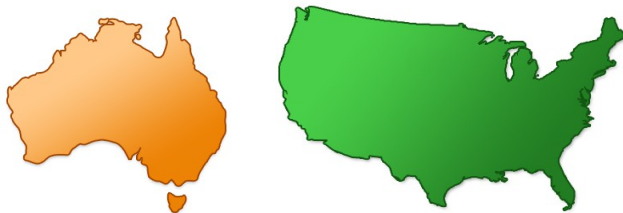
8.4 Handling Missing Data

Objectives

- Identify types of raw data files and input styles.
- Define the terms standard and nonstandard data.

Business Scenario

Information about Orion Star sales employees from Australia and the United States is stored in a raw data file.



Raw data file



Programmers need to be able to identify the layout and type of information in the raw data file.

Raw Data Files

A raw data file is also known as a *flat file*.

- They are text files that contain one record per line.
- A record typically contains multiple fields.
- Flat files do not have internal metadata.
- External documentation, known as a *record layout*, should exist.
- A record layout describes the fields and locations within each record.

Raw Data Files

Fields in a raw data file can be delimited or arranged in fixed columns.

Delimited File

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
```

Fixed Column File

1	1	2	2	3	3	4	4	5	5	6
1	5	0	5	0	5	0	5	0	5	0
120102	Tom		Zhou		Sales Manager		108255	AU		
120103	Wilson		Dawes		Sales Manager		87975	AU		
120121	Irenie		Elvish		Sales Rep. II		26600	AU		
120122	Christina		Ngan		Sales Rep. II		27475	AU		

Fields in Raw Data Files

In order for SAS to read a raw data file, you must specify the following information about each field:

- the location of the data value in the record
- the name of the SAS variable in which to store the data
- the type of the SAS variable

Reading Raw Data Files

There are different techniques, or *input styles*, for reading raw data files in SAS.

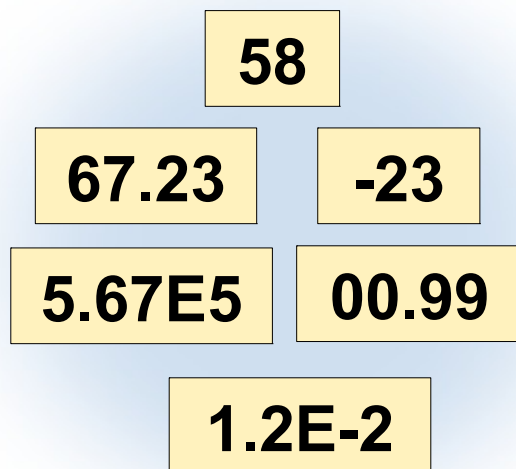
Input Style	Used for Reading
Column Input	Standard data in fixed columns
Formatted Input	Standard and nonstandard data in fixed columns
List Input	Standard and nonstandard data separated by blanks or some other delimiter

Standard and Nonstandard Data

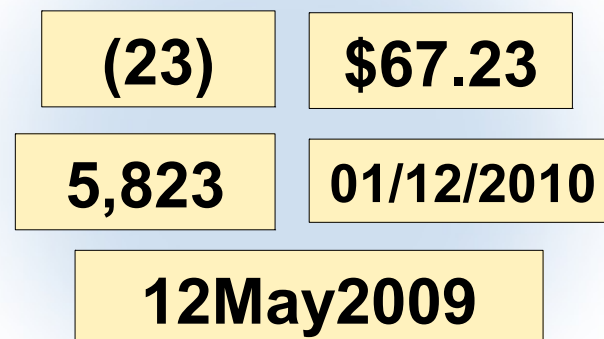
Standard data is data that SAS can read without any additional instruction.

- Character data is always standard.
- Some numeric values are standard and some are not.

Standard Numeric Data



Nonstandard Numeric Data



8.01 Multiple Answer Poll

What type of raw data files do you read?

- delimited
- fixed column
- both delimited and fixed column
- I do not read raw data files.

Chapter 8: Reading Raw Data Files

8.1 Introduction to Reading Raw Data Files

8.2 Reading Standard Delimited Data

8.3 Reading Nonstandard Delimited Data

8.4 Handling Missing Data

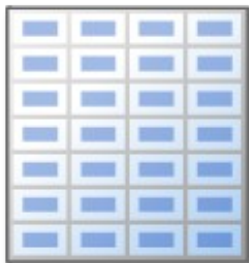
Objectives

- Use list input to create a SAS data set from a delimited raw data file.
- Examine the compilation and execution phases of the DATA step when reading a raw data file.
- Explicitly define the length of a variable.
- Examine behavior when a data error is encountered.

Business Scenario

Information about Orion Star sales employees is stored in a comma-delimited raw data file. The file contains both standard and nonstandard data fields.

sales.csv



DATA step



work.sales



List Input

Use list input to read delimited raw data files.

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
120123,Kimiko,Hotstone,F,26190,Sales Rep. I,AU,28SEP1968,10/01/1989
```

- SAS considers a space (blank) to be the default delimiter.
- Both standard and nonstandard data can be read.
- Fields must be read sequentially, left to right.

8.02 Quiz

Which fields in this file can be read as standard numeric values?

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
120123,Kimiko,Hotstone,F,26190,Sales Rep. I,AU,28SEP1968,10/01/1989
```

8.02 Quiz – Correct Answer

Which fields in this file can be read as standard numeric values?

The employee ID and salary. The date fields are nonstandard and require special processing.

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
120123,Kimiko,Hotstone,F,26190,Sales Rep. I,AU,28SEP1968,10/01/1989
```


Reading a Delimited Raw Data File

Use *INFILE* and *INPUT* statements in a DATA step to read a raw data file.

```
data work.subset;  
  infile "&path\sales.csv" dlm=',';  
  input Employee_ID First_Name $  
        Last_Name $ Gender $ Salary  
        Job_Title $ Country $;  
  
run;
```

```
DATA output-data-set;  
  INFILE "raw-data-file" <DLM='delimiter'>;  
  INPUT variable <$> variable <$> ... ;  
RUN;
```

INFILE Statement

The INFILE statement identifies the raw data file to be read.

```
INFILE "&path\sales.csv" DLM=' , ' ;
```

```
INFILE "raw-data-file" <DLM='delimiter'>;
```

- A full path is recommended.
- Using the **&path** macro variable reference makes the program more flexible.
- The DLM= option specifies alternate delimiters.



Be sure to use double quotation marks when referencing a macro variable within a quoted string.

INPUT Statement

The INPUT statement reads the data fields sequentially, left to right. Standard data fields require only a variable name and type.

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1969,06/01/1989
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1949,01/01/1974
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1944,01/01/1974
```

```
input Employee_ID First Name $ Last Name $
      Gender $ Salary Job Title $ Country $;
```

INPUT *variable* <\$> *variable* <\$> ...;

- The optional dollar sign indicates a character variable.
- Default length for ***all*** variables is eight bytes, regardless of type.

Viewing the Log

Partial SAS Log

```
249 data work.subset;  
250   infile "&path\sales.csv" dlm=',';  
251   input Employee_ID First_Name $ Last_Name $  
252         Gender $ Salary Job_Title $ Country $;  
253 run;
```

NOTE: The infile "s:\workshop\sales.csv" is:

Filename=s:\workshop\sales.csv,
RECFM=V,LRECL=256,File Size (bytes)=11340

NOTE: 165 records were read from the infile "s:\workshop\sales.csv".

The minimum record length was 61.

The maximum record length was 80.

NOTE: The data set WORK.SUBSET has 165 observations and 7 variables.

Viewing the Output

```
proc print data=work.subset noobs;
run;
```

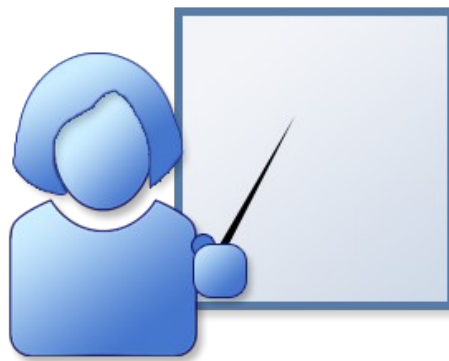
Partial PROC PRINT Output

Employee_ ID	First_ Name	Last_ Name	Gender	Salary	Job_ Title	Country
120102	Tom	Zhou	M	108255	Sales Ma	AU
120103	Wilson	Dawes	M	87975	Sales Ma	AU
120121	Irenie	Elvish	F	26600	Sales Re	AU
120122	Christin	Ngan	F	27475	Sales Re	AU
120123	Kimiko	Hotstone	F	26190	Sales Re	AU

Some character values are truncated.

Business Scenario

It is important to understand the processing that occurs when a DATA step reads a raw data file.



Compilation Phase

During compilation, SAS does the following:

- scans the step for syntax errors
- translates each statement into machine language
- creates an *input buffer* to hold one record at a time from the raw data file

Input Buffer										1									2
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6		8	9	0
																			...

- creates the program data vector (PDV) to hold one observation
- creates the descriptor portion of the output data set

Compilation

```
data work.subset;  
    infile "&path\sales.csv" dlm=',';  
    input Employee_ID First_Name $ Last_Name $  
          Gender $ Salary Job_Title $ Country $;  
run;
```


Compilation

```
data work.subset;
  infile "&path\sales.csv" dlm=',';
  input Employee_ID First Name $ Last Name $
        Gender $ Salary Job_Title $ Country $;
run;
```

Input Buffer

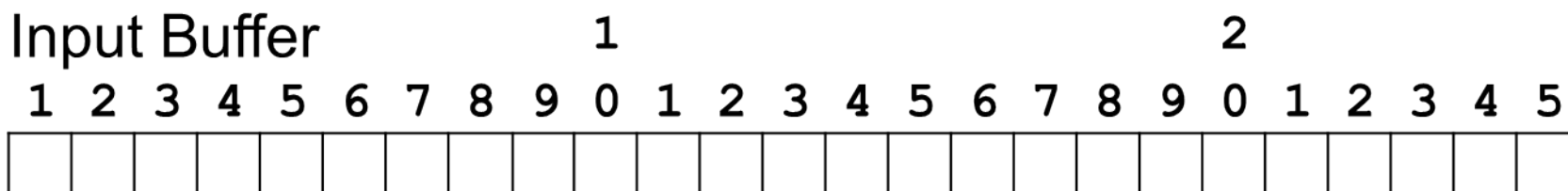
1

2

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5

Compilation

```
data work.subset;
  infile "&path\sales.csv" dlm=',';
  input Employee_ID First Name $ Last Name $
        Gender $ Salary Job_Title $ Country $;
run;
```



PDV

Employee_ID
_ID
N 8

Attributes are based on the INPUT statement.

Compilation

```
data work.subset;
  infile "&path\sales.csv" dlm=',';
  input Employee_ID First_Name $ Last_Name $
         Gender $ Salary Job_Title $ Country $;
run;
```

Input Buffer

Input Buffer

1										2														
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5

PDV

Employee_ID N 8	First_Name \$ 8

With list input, the default length for character variables is eight bytes.

Compilation

```
data work.subset;
  infile "&path\sales.csv" dlm=',';
  input Employee ID First Name $ Last Name $
        Gender $ Salary Job_Title $ Country $;
run;
```

Input Buffer

Input Buffer

1

2

1

2

3

4

5

6

7

8

9

0

1

2

3

4

5

6

7

8

9

0

1

2

3

4

5

PDV

Employee _ID N 8	First_ Name \$ 8	Last _Name \$ 8	Gender \$ 8	Salary N 8	Job_ _Title \$ 8	Country \$ 8

Compilation

```
data work.subset;
  infile "&path\sales.csv" dlm=',';
  input Employee_ID First Name $ Last Name $
        Gender $ Salary Job_Title $ Country $;
run;
```

PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8

Descriptor Portion of **work.subset**

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8

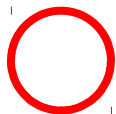
8.03 Multiple Choice Poll

Which statement is true?

- An input buffer is created only if you are reading data from a raw data file.
- The PDV at compile time holds the variable name, type, byte size, and initial value.
- The descriptor portion is the first item that is created at compile time.

8.03 Multiple Choice Poll – Correct Answer

Which statement is true?



- An input buffer is created only if you are reading data from a raw data file.
- The PDV at compile time holds the variable name, type, byte size, and initial value.
- The descriptor portion is the first item that is created at compile time.

DATA Step Processing

Compile the step

Success?

No

Next step

Yes

Compilation Phase

Execution Phase

Initialize all variables to missing

Execute INPUT statement

Execute other statements

Output to SAS data set

End of file?

Yes

Next step

No

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

Initialize PDV

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

Input Buffer

Input Buffer

										1											2						
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5			

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
    dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

Input Buffer

Input Buffer

1

2

1

2

3

4

5

6

7

8

9

0

1

2

3

4

5

6

7

8

9

0

1

2

3

4

5

PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8
.				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

run:

SAS reads a record into the input buffer.

Input Buffer

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
.				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
    dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

SAS scans until it reaches a delimiter.

Input Buffer

1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5

1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
.				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job Title $
```

run

The value is converted from text to a floating-point numeric value and copied to the PDV.

Input Buffer

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job Title \$ 8	Country \$ 8
120102				.		

Execution

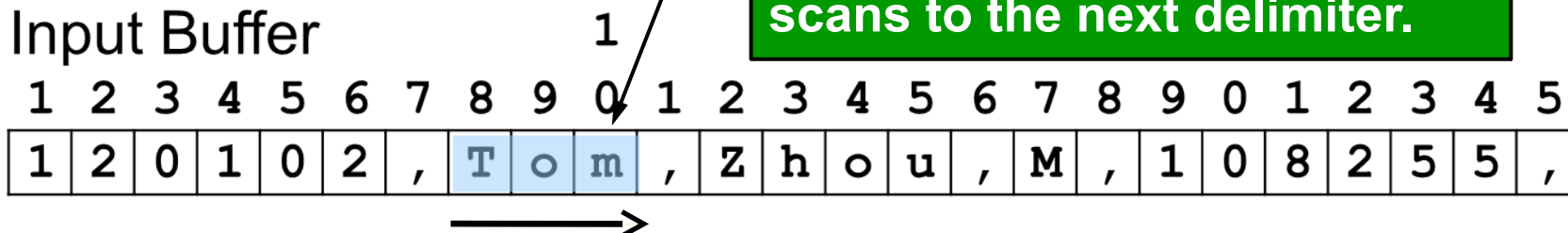
Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
    dlm=',';
  input Employee_ID First_Name $
    Last_Name $ Gender $
    Salary Job_Title $
    Country $;
```

SAS skips the delimiter and scans to the next delimiter.

Input Buffer



PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8
120102				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
```

run;

The text value is copied to the PDV without conversion.

Input Buffer

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8
120102	Tom			.		

Partial sales.csv

```
data work.subset;  
    infile "&path\sales.csv"  
        dlm=',';  
    input Employee_ID First_Name $  
        Last_Name $ Gender $  
        Salary Job_Title $  
        Country $;  
run;
```

```
run ;
```

Input Buffer

1

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
120102	Tom	Zhou	M	108255	Sales Ma	AU

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

**Implicit OUTPUT;
Implicit RETURN;**

Input Buffer

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
120102	Tom	Zhou	M	108255	Sales Ma	AU

Execution

Here is the output data set after the first iteration of the DATA step.

work.subset

Employee _ID	First _Name	Last _Name	Gende r	Salary	Job _Title	Country
120102	Tom	Zhou	M	108255	Sales Ma	AU

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

**Implicit OUTPUT;
Implicit RETURN;**

Input Buffer

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8
120102	Tom	Zhou	M	108255	Sales Ma	AU

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

Reinitialize PDV

```
data work.subset;
  infile "&path\sales.csv"
    dlm=',';
  input Employee_ID First_Name $
    Last_Name $ Gender $
    Salary Job_Title $
    Country $;

run;
```

Input Buffer

Input Buffer

1										2														
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

All variables in the PDV are reinitialized.

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8
.				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

Input Buffer

Input Buffer

1										2														
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	2	,	T	o	m	,	Z	h	o	u	,	M	,	1	0	8	2	5	5	,

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
.				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;

run;
```

Input Buffer

Input Buffer

1										2														
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	3	,	W	i	l	s	o	n	,	D	a	w	e	s	,	M	,	8	7	9

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
.				.		

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

Input Buffer

Input Buffer

1										2														
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	3	,	W	i	l	s	o	n	,	D	a	w	e	s	,	M	,	8	7	9

PDV

Employee_ID N 8	First_Name \$ 8	Last_Name \$ 8	Gender \$ 8	Salary N 8	Job_Title \$ 8	Country \$ 8
120103	Wilson	Dawes	M	87975	Sales Ma	AU

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "&path\sales.csv"
        dlm=',';
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

**Implicit OUTPUT;
Implicit RETURN;**

Input Buffer

1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	3	,	W	i	l	s	o	n	,	D	a	w	e	s	,	M	,	8	7	9

work.subset

Employee_ID	First_Name	Last_Name	Gender	Salary	Job_Title	Country
120102	Tom	Zhou	M	108255	Sales Ma	AU
120103	Wilson	Dawes	M	87975	Sales Ma	AU

Execution

Partial sales.csv

```
120102,Tom,Zhou, ...
120103,Wilson,Dawes, ...
120121,Irenie,Elvish, ...
120122,Christina,Ngan, ...
120123,Kimiko,Hotstone, ...
120124,Lucian,Daymond, ...
120125,Fong,Hofmeister, ...
```

```
data work.subset;
  infile "smith\sales.csv"
  input Employee_ID First_Name $
        Last_Name $ Gender $
        Salary Job_Title $
        Country $;
run;
```

Continue until EOF

Input Buffer										1					2									
1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5
1	2	0	1	0	3	,	W	i	l	s	o	n	,	D	a	w	e	s	,	M	,	8	7	9

PDV

Employee _ID N 8	First _Name \$ 8	Last _Name \$ 8	Gende r \$ 8	Salary N 8	Job _Title \$ 8	Country \$ 8
.				.		

8.04 Multiple Choice Poll

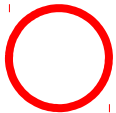
Which statement is true of a DATA step when reading from a raw data file?

- Data is read from the raw data file into the PDV.
- The size of the input buffer adjusts automatically based on the length of the input record.
- At the bottom of the DATA step, the contents of the PDV are output to the output SAS data set.

8.04 Multiple Choice Poll – Correct Answer

Which statement is true of a DATA step when reading from a raw data file?

- Data is read from the raw data file into the PDV.
- The size of the input buffer adjusts automatically based on the length of the input record.
- At the bottom of the DATA step, the contents of the PDV are output to the output SAS data set.



Viewing the Output

```
proc print data=work.subset noobs;  
  
run;
```

Partial PROC PRINT Output

Employee_ ID	First_ Name	Last_ Name	Gender	Salary	Job_ Title	Country
120102	Tom	Zhou	M	108255	Sales Ma	AU
120103	Wilson	Dawes	M	87975	Sales Ma	AU
120121	Irenie	Elvish	F	26600	Sales Re	AU
120122	Christin	Ngan	F	27475	Sales Re	AU
120123	Kimiko	Hotstone	F	26190	Sales Re	AU
120124	Lucian	Daymond	M	26480	Sales Re	AU
120125	Fong	Hofmeister	M	32040	Sales Re	AU

Some character values are truncated.

LENGTH Statement

The *LENGTH* statement defines the type and length of a variable.

```
data work.subset;  
    length First_Name $ 12 Last_Name $ 18  
           Gender $ 1 Job_Title $ 25  
           Country $ 2;  
    infile "&path\sales.csv" dlm=',';  
    input Employee_ID First_Name $ Last_Name $  
           Gender $ Salary Job_Title $ Country $;  
run;
```

LENGTH *variable(s) \$ length;*



Put the LENGTH statement before the INPUT statement.

Compilation

```
data work.subset;
  length First_Name $ 12 Last_Name $ 18
         Gender $ 1 Job_Title $ 25
         Country $ 2;
  infile "&path\sales.csv" dlm=',';
  input Employee_ID First_Name $ Last_Name $
         Gender $ Salary Job_Title $ Country $;
run;
```

PDV

Attributes are based on the LENGTH statement.

First_ Name \$ 12	Last _Name \$ 18	Gender \$ 1	Job_Title \$ 25	Country \$ 2

Compilation

```
data work.subset;
  length First_Name $ 12 Last_Name $ 18
         Gender $ 1 Job_Title $ 25
         Country $ 2;
  infile "&path\sales.csv" dlm=',';
  input Employee_ID First_Name $ Last_Name $
        Gender $ Salary Job_Title $ Country $;
run;
```

PDV

First_ Name \$ 12	Last _Name \$ 18	Gender \$ 1	Job_Title \$ 25	Country \$ 2	Employee_ ID N 8	Salary N 8



Viewing the Output

```
proc print data=work.subset noobs;  
run;
```

Partial PROC PRINT Output

First_	Employee_					
Name	Last_Name	Gender	Job_Title	Country	ID	Salary
Tom	Zhou	M	Sales Manager	AU	120102	108255
Wilson	Dawes	M	Sales Manager	AU	120103	87975
Irenie	Elvish	F	Sales Rep. II	AU	120121	26600
Christina	Ngan	F	Sales Rep. II	AU	120122	27475
Kimiko	Hotstone	F	Sales Rep. I	AU	120123	26190

The character values are no longer truncated, but the order of the variables has changed.

8.05 Quiz

Suppose you want the order of the variables to match the order of the fields. You can include the numeric variables in the LENGTH statement. Which of the following produces the correct results?

a. `length Employee_ID First_Name $ 12
Last_Name $ 18 Gender $ 1
Salary Job_Title $ 25
Country $ 2;`

b. `length Employee_ID 8 First_Name $ 12
Last_Name $ 18 Gender $ 1
Salary 8 Job_Title $ 25
Country $ 2;`

8.05 Quiz – Correct Answer

Suppose you want the order of the variables to match the order of the fields. You can include the numeric variables in the LENGTH statement. Which of the following produces the correct results?

a. `length Employee_ID First_Name $ 12
Last_Name $ 18 Gender $ 1
Salary Job_Title $ 25
Country $ 2;`

b. `length Employee_ID 8 First_Name $ 12
Last_Name $ 18 Gender $ 1
Salary 8 Job_Title $ 25
Country $ 2;`

Using a LENGTH Statement

The LENGTH statement identifies the character variables, so dollar signs can be omitted from the INPUT statement.

```
data work.subset;  
    length Employee_ID 8 First_Name $ 12  
           Last_Name $ 18 Gender $ 1  
           Salary 8 Job_Title $ 25  
           Country $ 2;  
    infile "&path\sales.csv" dlm=',';  
    input Employee_ID First_Name Last_Name  
          Gender Salary Job_Title Country;  
run;
```

Viewing the Output

Display the variables in creation order.

```
proc contents data=work.subset varnum;  
run;
```

Partial PROC CONTENTS Output

Variables in Creation Order

#	Variable	Type	Len
1	Employee_ID	Num	8
2	First_Name	Char	12
3	Last_Name	Char	18
4	Gender	Char	1
5	Salary	Num	8
6	Job_Title	Char	25
7	Country	Char	2

Viewing the Output

Partial PROC PRINT Output

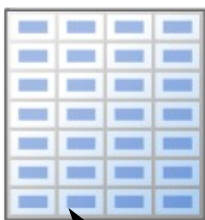
Employee_	First_					
ID	Name	Last_Name	Gender	Salary	Job_Title	Country
120102	Tom	Zhou	M	108255	Sales Manager	AU
120103	Wilson	Dawes	M	87975	Sales Manager	AU
120121	Irenie	Elvish	F	26600	Sales Rep. II	AU
120122	Christina	Ngan	F	27475	Sales Rep. II	AU
120123	Kimiko	Hotstone	F	26190	Sales Rep. I	AU
120124	Lucian	Daymond	M	26480	Sales Rep. I	AU
120125	Fong	Hofmeister	M	32040	Sales Rep. IV	AU



Business Scenario

A raw data file contains information about Orion Star sales employees. It includes some invalid data values.

sales3inv.csv

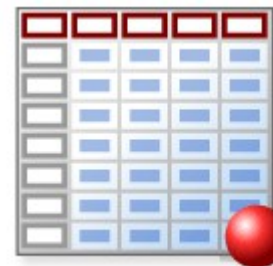


invalid values

DATA step



work.sales



8.06 Quiz

What problems do you see with the data values for the last two data fields, **Salary** and **Country**?

Partial **sales3inv.csv**

```
120102,Tom,Zhou,Manager,108255,AU
120103,Wilson,Dawes,Manager,87975,AU
120121,Irenie,Elvish,Rep. II,26600,AU
120122,Christina,Ngan,Rep. II,n/a,AU
120123,Kimiko,Hotstone,Rep. I,26190,AU
120124,Lucian,Daymond,Rep. I,26480,12
120125,Fong,Hofmeister,Rep. IV,32040,AU
```


8.06 Quiz – Correct Answer

What problems do you see with the data values for the last two data fields, **Salary** and **Country**?

Partial **sales3inv.csv**

120102	Tom	Zhou	Manager	108255	AU
120103	Wilson	Dawes	Manager	87975	AU
120121	Irenie	Elvish	Rep. II	26600	AU
120122	Christina	Ngan	Rep. II	n/a	AU
120123	Kimiko	Hotstone	Rep. I	26190	AU
120124	Lucian	Daymond	Rep. I	26480	12
120125	Fong	Hofmeister	Rep. IV	32040	AU

Reading a Raw Data File with Data Errors

```
data work.sales;  
    infile "&path\sales3inv.csv" dlm=',';  
    input Employee_ID First $ Last $  
           Job_Title $ Salary Country $;  
run;  
  
proc print data=work.sales;  
run;
```

Salary is defined as numeric and **Country** as character.

Viewing the Output

Partial PROC PRINT Output

Obs	Employee_ ID	First	Last	Job_ Title	Salary	Country
1	120102	Tom	Zhou	Manager	108255	AU
2	120103	Wilson	Dawes	Manager	87975	AU
3	120121	Irenie	Elvish	Rep. II	26600	AU
4	120122	Christina	Ngan	Rep. II	.	AU
5	120123	Kimiko	Hotstone	Rep. I	26190	AU
6	120124	Lucian	Daymond	Rep. I	26480	12
7	120125	Fong	Hofmeister	Rep. IV	32040	AU

- A missing value was stored in **Salary** for the input value *n/a*.
- The value *12* was successfully stored in **Country**.
- A data error occurred on observation 4 but not on observation 6.

Viewing the Log

Partial SAS Log

```
480 data work.sales;  
481   infile "&path\sales3inv.csv" dlm=',';  
482   input Employee_ID First $ Last $  
483         Job_Title $ Salary Country $;  
484 run;
```

NOTE: The infile "s:\workshop\sales3inv.csv" is:
 Filename=s:\workshop\sales3inv.csv,
 RECFM=V,LRECL=256,File Size (bytes)=1972,

NOTE: Invalid data for Salary in line 4 31-33.

RULE: ----+----1----+----2----+----3----+----4----+----5-

4 120122,Christina,Ngan,Rep. II,n/a,AU 36

Employee_ID=120122 First=Christin Last=Ngan Job_Title=Rep. II Salary=
Country=AU _ERROR_=1 _N_=4

NOTE: 50 records were read from the infile "s:\workshop\sales3inv.csv".

NOTE: The data set WORK.SALES has 50 observations and 6 variables.

A data error occurs when a data value does not match the field specification.

Data Errors

When this kind of data error occurs, the following information is written to the SAS log:

- a note describing the error
- a column ruler
- the input record
- the contents of the PDV

NOTE: Invalid data for Salary in line 4 31-33.

RULE: ----+----1----+----2----+----3----+----4----+----5-

4 120122,Christina,Ngan,Rep. II,n/a,AU 36

Employee_ID=120122 First=Christin Last=Ngan Job_Title=Rep. II Salary=.

Country=AU _ERROR_=1 _N_=4

A missing value is assigned to the corresponding variable, and execution continues.

Data Errors

Two temporary variables are created during the processing of every DATA step:

- **_N_** is the DATA step iteration counter.
- **_ERROR_** indicates data error status.
 - 0 indicates that no data error occurred on that record.
 - 1 indicates that one or more data errors occurred on that record.

NOTE: Invalid data for Salary in line 4 31-33.

RULE: ----+----1----+----2----+----3----+----4----+----5-

4 120122,Christina,Ngan,Rep. II,n/a,AU 36

Employee_ID=120122 First=Christin Last=Ngan Job_Title=Rep. II Salary=.

Country=AU **_ERROR_=1 _N_=4**



Examining Data Errors

This demonstration illustrates SAS behavior when a data error occurs when reading a raw data file.

8.07 Multiple Choice Poll

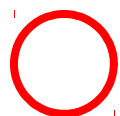
Submit program **p108a01** and examine the log.

Which statement best describes the reason for the error?

- The data in the raw data file is invalid.
- The programmer incorrectly read the data.

8.07 Multiple Choice Poll – Correct Answer

Which statement best describes the reason for the error?



- The data in the raw data file is invalid.
- The programmer incorrectly read the data.

Partial SAS Log

```
404      input Employee_ID First $ Last;
405      run;
```

Last was read as
numeric but needs to be
read as character.

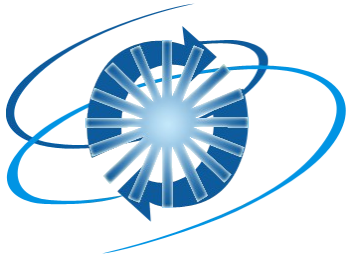
NOTE: Invalid data for Last in line 1 16-17.

```
RULE:      -----1-----2-----3-----4-----5-----6
1          120101,Patrick,Lu,M,163040,Director,AU,18AUG1976,01JUL2003 58
Employee_ID=120101 First=Patrick Last=._ERROR_=1 _N_=1
```

NOTE: Invalid data for Last in line 2 15-24.

```
2          120104,Kareen,Billington,F,46230,Administration Manager,au,1
61 1MAY1954,01JAN1981 78
Employee_ID=120104 First=Kareen Last=._ERROR_=1 _N_=2
```





Exercise

This exercise reinforces the concepts discussed previously.

Chapter 8: Reading Raw Data Files

8.1 Introduction to Reading Raw Data Files

8.2 Reading Standard Delimited Data

8.3 Reading Nonstandard Delimited Data

8.4 Handling Missing Data

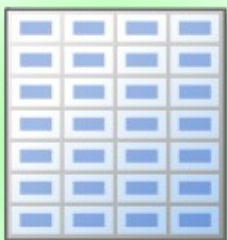
Objectives

- Use informats to read character data.
- Use informats to read nonstandard data.
- Subset observations and add permanent attributes.

Business Scenario

Create a temporary SAS data set by reading both standard and nonstandard values from a comma-delimited raw data file.

sales.csv



DATA step



work.sales



The new data set will contain a subset of the input data and will include permanent attributes.

Considerations

Use modified list input to read all the fields from **sales.csv**.
Store the date fields as SAS dates.

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
120123,Kimiko,Hotstone,F,26190,Sales Rep. I,AU,28SEP1968,10/01/1989
```

Modified List Input

This DATA step uses *modified list input*. Instead of a LENGTH statement, an informat specifies the length for each character variable.

```
data work.subset;  
  infile "&path\sales.csv" dlm=',';  
  input Employee_ID First Name :$12.  
         Last Name :$18. Gender :$1. Salary  
         Job Title :$25. Country :$2.;  
run;
```

input variable <:informat.> ...;

- The **\$12.** informat defines a length of 12 for **First Name** and allows up to 12 characters to be read.
- The **:** format modifier tells SAS to read until it encounters a delimiter.

Modified List Input



Omitting the colon modifier causes unexpected results.

Partial **sales.csv**

reads 12 characters

120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993

```
input Employee_ID First_Name $12.
      Last_Name   :$18. Gender   :$1. Salary
      Job_Title   :$25. Country  :$2.;
```

PDV

Employee_ID N 8	First_Name \$ 12	Last_Name \$ 18	Gender \$ 1
120102	Tom,Zhou,1	08255	S

Salary N 8	Job_Title \$ 25	Country \$ 2
.	11AUG1973	06

Reading Nonstandard Data

An informat is ***required*** to read nonstandard numeric data.

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
120123,Kimiko,Hotstone,F,26190,Sales Rep. I,AU,28SEP1968,10/01/1989
```

In this example, informats are needed to specify the style of the date fields so that they can be read and converted to SAS dates.

8.08 Quiz

A ***format*** is an instruction that tells SAS how to display data values. What formats would you specify to display a SAS date in the styles shown below?

a) 01JAN2000

b) 01/16/2000

8.08 Quiz – Correct Answer

A ***format*** is an instruction that tells SAS how to display data values. What formats would you specify to display a SAS date in the styles shown below?

a) 01JAN2000 ☐ **DATE9.**

b) 01/16/2000 ☐ **MMDDYY10.**

Informats for Nonstandard Data

An *informat* is an instruction that SAS uses to **read** data values into a variable.

Partial **sales.csv**

```
120102,Tom,Zhou,M,108255,Sales Manager,AU,11AUG1973,06/01/1993
120103,Wilson,Dawes,M,87975,Sales Manager,AU,22JAN1953,01/01/1978
120121,Irenie,Elvish,F,26600,Sales Rep. II,AU,02AUG1948,01/01/1978
120122,Christina,Ngan,F,27475,Sales Rep. II,AU,27JUL1958,07/01/1982
120123,Kimiko,Hotstone,F,26190,Sales Rep. I,AU,28SEP1968,10/01/1989
```

DATE.

MMDDYY.

The informat describes the data value and tells SAS how to convert it

SAS Informats

SAS informats have the following form:

```
<$><informat><w>.
```

\$	Indicates a character informat.
<i>informat</i>	Names the SAS informat or user-defined informat.
<i>w</i>	Specifies the width or number of columns to read or specifies the length of a character variable.
.	Is required syntax.

- The width is typically not used with list input because SAS will read each field until it encounters a delimiter.

SAS Informats

Selected SAS Informats for Nonstandard Numeric Values

Informat	Definition
COMMA. DOLLAR.	Reads nonstandard numeric data and removes embedded commas, blanks, dollar signs, percent signs, and dashes.
COMMAX. DOLLARX.	Reads nonstandard numeric data and removes embedded non-numeric characters; reverses the roles of the decimal point and the comma.
EUROX.	Reads nonstandard numeric data and removes embedded non-numeric characters in European currency.
\$CHAR.	Reads character values and preserves leading blanks.
\$UPCASE.	Reads character values and converts them to uppercase.

SAS Informats

Informats are used to read and convert raw data.

Informat	Raw Data Value	SAS Data Value
COMMA. DOLLAR.	\$12,345	12345
COMMAX. DOLLARX.	\$12.345	12345
EUROX.	€12.345	12345
\$CHAR.	##Australia	##Australia
\$UPCASE.	au	AU

The character # represents a blank space.

SAS Informats

Use date informats to read and convert dates to SAS date values.

Informat	Raw Data Value	SAS Data Value
MMDDYY.	010160 01/01/60 01/01/1960 1/1/1960	0
DDMMYY.	311260 31/12/60 31/12/1960	365
DATE.	31DEC59 31DEC1959	-1

8.09 Quiz

Use the SAS Help facility or documentation to investigate the **DATEw.** informat and answer the following questions:

- a) What does the **w** represent?

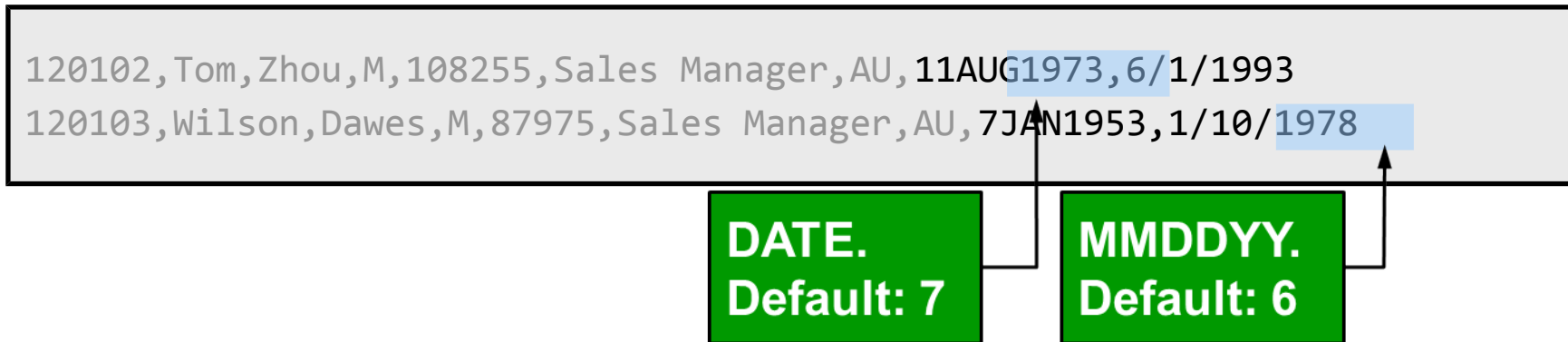
- b) What is the default width of this informat?

8.09 Quiz – Correct Answer

Use the SAS Help facility or documentation to investigate the **DATEw.** informat and answer the following questions:

- a) What does the **w** represent?
the width of the input field
- b) What is the default width of this informat?
The default width is 7.

Using Informats to Read Nonstandard Data



- An informat is needed to read a nonstandard value.
- There is no need to specify a width with list input.
- Most informats have default widths.



Only the highlighted portion of the field is read if the default width is specified.

Modified List Input

The colon format modifier (:) tells SAS to read until it encounters a delimiter.

```
input Employee_ID First_Name :$12.  
      Last_Name :$18. Gender :$1.  
      Salary Job_Title :$25. Country :$2.  
      Birth_Date :date. Hire Date :mmddyy.;
```

INPUT *variable* <\$> *variable* <:informat> ...;

colon format modifier

Viewing the Log

```
37 data work.sales;  
38   infile "&path\sales.csv" dlm=',';  
39   input Employee_ID First_Name :$12. Last_Name :$18.  
40         Gender :$1. Salary Job_Title :$25. Country :$2.  
41         Birth_Date :date. Hire_Date :mmddyy.;  
42 run;
```

NOTE: The infile "s:\workshop\sales.csv" is:

**Filename=s:\workshop\sales.csv,
RECFM=V,LRECL=256,File Size (bytes)=11340,**

NOTE: 165 records were read from the infile "s:\workshop\sales.csv".

NOTE: The data set WORK.SALES has 165 observations and 9 variables.

Viewing the Output

```
proc print data=work.sales;
run;
```

Partial PROC PRINT Output

Obs	First_ Name	Last_Name	Gender	Job_Title	Country	Employee_ ID	Salary	Birth_ Date	Hire_ Date
1	Tom	Zhou	M	Sales Manager	AU	120102	108255	4971	12205
2	Wilson	Dawes	M	Sales Manager	AU	120103	87975	-2535	6575
3	Irenie	Elvish	F	Sales Rep. II	AU	120121	26600	-4169	6575
4	Christina	Ngan	F	Sales Rep. II	AU	120122	27475	-523	8217
5	Kimiko	Hotstone	F	Sales Rep. I	AU	120123	26190	3193	10866

8.10 Multiple Choice Poll

A new data set should contain observations only for Australian employees. Which of the following can be used to subset the data?

- `where Country='AU';`
- `if Country='AU';`
- either a or b
- You cannot subset when reading from a raw data file.

8.10 Multiple Choice Poll – Correct Answer

A new data set should contain observations only for Australian employees. Which of the following can be used to subset the data?

- ☒ – where Country='AU';
- ☐ – if Country='AU';
- ☐ – either a or b
- ☐ – You cannot subset when reading from a raw data file.

Additional SAS Statements

Additional SAS statements can be added to perform further processing in the DATA step.

```
data work.sales;  
  infile "&path\sales.csv" dlm=',';  
  input Employee_ID First_Name :$12. Last_Name :$18.  
        Gender :$1. Salary Job_Title :$25. Country :$2.  
        Birth_Date :date. Hire_Date :mmddy.;  
  if Country='AU';  
  keep First_Name Last_Name Salary  
        Job_Title Hire_Date;  
  label Job_Title='Sales Title'  
        Hire_Date='Date Hired';  
  format Salary dollar12. Hire_Date monyy7.;  
run;
```


Viewing the Output

```
proc print data=work.sales label;
run;
```

Partial PROC PRINT Output

	First_ Salary	Hired		Date Obs	Name	Last_Name	Sales Title
1	Tom	Zhou	Sales Manager	\$108,255	JUN1993		
2	Wilson	Dawes	Sales Manager	\$87,975	JAN1978		
3	Irenie	Elvish	Sales Rep. II	\$26,600	JAN1978		
4	Christina	Ngan	Sales Rep. II	\$27,475	JUL1982		
5	Kimiko	Hotstone	Sales Rep. I	\$26,190	OCT1989		

WHERE versus Subsetting IF Statement

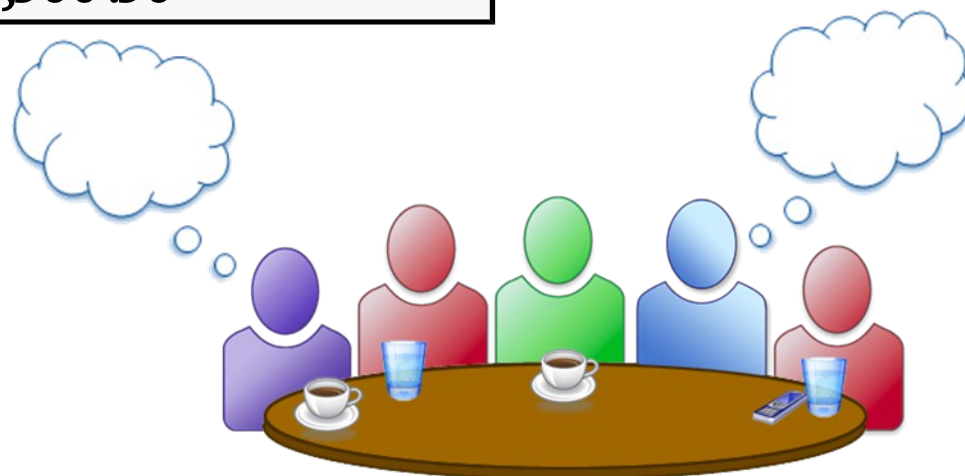
Step and Usage	WHERE	IF
PROC step	Yes	No
DATA step (source of variable)		
SET statement	Yes	Yes
assignment statement	No	Yes
 INPUT statement	No	Yes

Idea Exchange

What factors need to be considered when reading **salary.dat**, shown below?

Partial **salary.dat**

	1	1	2	2	3	3
	1	5	0	5	0	5
Donny	5	M	AY	2008	25	FL \$43,132.50
Margaret	20	FEB	2008	43	NC	\$65,150
Dan	1	JUN	2008	27	FL	\$40,000.00
Subash	2	FEB	2008	45	NC	\$75,750
Antonio	25	M	AY	2008	35	FL \$43,500.50





Using List Input: Importance of Colon Format Modifier

This demonstration illustrates the use of informats when reading nonstandard data values.

Business Scenario

You are working on a new project, but the raw data file has not been created yet. You can include in-stream data in a DATA step.



DATALINES Statement

The DATALINES statement supplies data within a program.

```
data work.newemps;  
    input First_Name $ Last_Name $  
          Job_Title $ Salary :dollar8.;  
datalines;  
Steven Worton Auditor $40,450  
Merle Hieds Trainee $24,025  
Marta Bamberger Manager $32,000  
;
```

DATALINES;
...
;

- DATALINES is the last statement in the DATA step and immediately precedes the first data line.
- A null statement (a single semicolon) indicates the end of the input data.

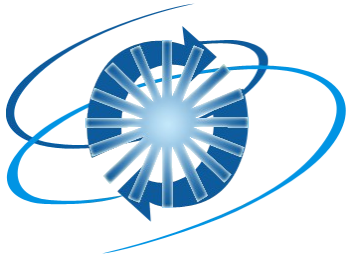
Viewing the Output

```
proc print data=work.newemps;  
run;
```

PROC PRINT Output

	First_	Last_	Job_	
Obs	Name	Name	Title	Salary
1	Steven	Worton	Auditor	40450
2	Merle	Hieds	Trainee	24025
3	Marta	Bamberge	Manager	32000





Exercise

This exercise reinforces the concepts discussed previously.

Chapter 8: Reading Raw Data Files

8.1 Introduction to Reading Raw Data Files

8.2 Reading Standard Delimited Data

8.3 Reading Nonstandard Delimited Data

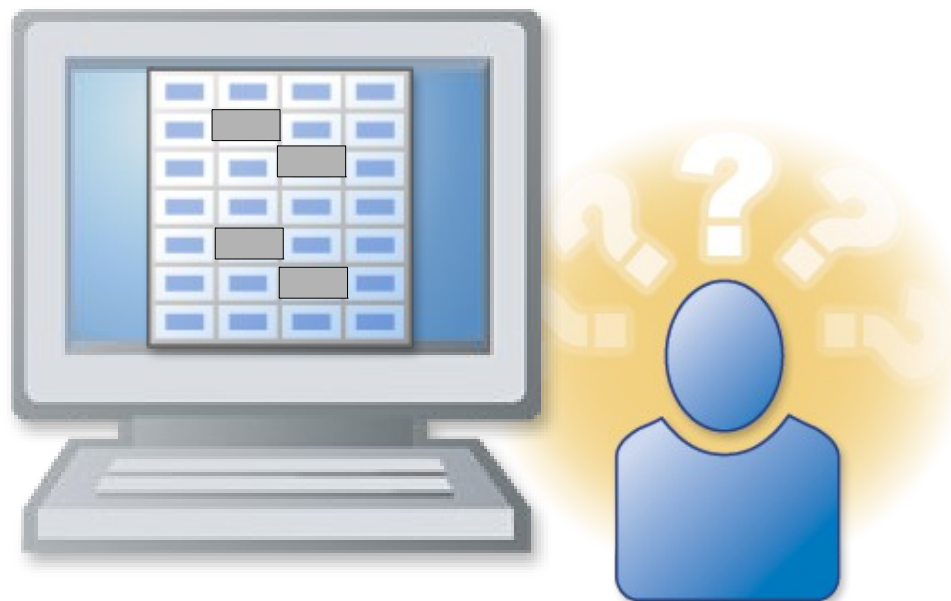
8.4 Handling Missing Data

Objectives

- Use the DSD option to read consecutive delimiters as missing values.
- Use the MISSOVER option to recognize missing values at the end of a record.

Business Scenario

Orion Star programmers have discovered that some files have records with missing data in one or more fields.



Missing Values in the Middle of the Record

The records in **phone2.csv** have a contact name, phone number, and a mobile number. The phone number is missing from some of the records.

Missing data is indicated by consecutive delimiters.

phone2.csv

	1	1	2	2	3	3	4	4		
1	---	5	---	0	---	5	---	0	---	5
James Kvarniq,	(704)	293-8126,	(701)	281-8923						
Sandrina Stephano,	(019)	271-4592								
Cornelia Krah1,	(212)	891-3241,	(212)	233-5413						
Karen Ballinger,	(714)	644-9090								
Elke Wallstab,	(910)	763-5561,	(910)	545-3421						

8.11 Quiz

- Open and submit **p108a03**.
- Examine the SAS log. How many input records were read and how many observations were created?
- Examine the report. Does it look correct?

```
data work.contacts;  
  length Name $ 20 Phone Mobile $ 14;  
  infile "&path\phone2.csv" dlm=',';  
  input Name $ Phone $ Mobile $;  
run;  
  
proc print data=work.contacts noobs;  
run;
```


8.11 Quiz – Correct Answer

- How many input records were read and how many observations were created? **five read, three created**
- Examine the report. Does it look correct? **no**

NOTE: 5 records were read from the infile "S:\workshop\phone2.csv".
 The minimum record length was 31.
 The maximum record length was 44.

NOTE: SAS went to a new line when INPUT statement reached past the end of a line.

NOTE: The data set WORK.CONTACTS has 3 observations and 3 variables.

Name	Phone	Mobile
James Kvarniq	(704) 293-8126	(701) 281-8923
Sandrina Stephano	(919) 871-7830	Cornelia Krah1
Karen Ballinger	(714) 344-4321	Elke Wallstab

Consecutive Delimiters in List Input

List input treats two or more consecutive delimiters as a single delimiter and not as a missing value.

phone2.csv

	1	1	2	2	3	3	4	4		
1	---	5	---	0	---	5	---	0	---	5
James Kvarniq,	(704)	293-8126,	(701)	281-8923						
Sandrina Stephano,	(919)	271-4592								
Cornelia Krah1,	(212)	891-3241,	(212)	233-5413						
Karen Ballinger,	(714)	644-9090								
Elke Wallstab,	(910)	763-5561,	(910)	545-3421						

When there is missing data in a record, SAS does the following:

- loads the next record to finish the observation
- writes a note to the log

DSD Option

Use the DSD option to correctly read **phone2.csv**.

```
data work.contacts;  
  length Name $ 20 Phone Mobile $ 14;  
  infile "&path\phone2.csv" dsd;  
  input Name $ Phone $ Mobile $;  
run;  
  
proc print data=work.contacts noobs;  
run;
```

```
INFILE "raw-data-file" <DLM=> DSD;
```

DSD Option

The DSD option in the INFILE statement does the following:

- sets the default delimiter to a comma
- treats consecutive delimiters as missing values
- enables SAS to read values with embedded delimiters if the value is surrounded by quotation marks

```
INFILE "raw-data-file" <DLM=> DSD;
```

The DLM= option can be used with the DSD option but is not needed for comma-delimited files.

Viewing the Output

Adding the DSD option gives the correct results.

PROC PRINT Output

Name	Phone	Mobile
James Kvarniq	(704) 293-8126	(701) 281-8923
Sandrina Stephano		(919) 271-4592
Cornelia Krah1	(212) 891-3241	(212) 233-5413
Karen Ballinger		(714) 644-9090
Elke Wallstab	(910) 763-5561	(910) 545-3421

Partial SAS Log

NOTE: 5 records were read from the infile "S:\workshop\phone2.csv".
The minimum record length was 31.
The maximum record length was 44.
NOTE: The data set WORK.CONTACTS has 5 observations and 3 variables.



Business Scenario

Orion Star programmers have discovered that some files have observations with missing data at the end of the record, so there are fewer fields in the record than specified in the INPUT statement.



Missing Values at the End of a Record

The raw data file **phone.csv** contains missing values at the end of some records.

phone.csv

1	1	2	2	missing values				4		
1	---	5	---	0	---	5	---	0	---	5
James Kvarniq,(704) 293-8126,(701) 281-8923										
Sandrina Stephano,(919) 871-7830										
Cornelia Krah1,(212) 891-3241,(212) 233-5413										
Karen Ballinger,(714) 344-4321										
Elke Wallstab,(910) 763-5561,(910) 545-3421										

The DSD option is not appropriate because the missing data is not marked by consecutive delimiters.

MISSOVER Option

The *MISSOVER* option prevents SAS from loading a new record when the end of the current record is reached.

```
data contacts;  
  length Name $ 20 Phone Mobile $ 14;  
  infile "&path\phone.csv" dlm=', ' missover;  
  input Name $ Phone $ Mobile $;  
run;  
  
proc print data=contacts noobs;  
run;
```

```
INFILE "raw-data-file" <DLM=> MISSOVER;
```

If SAS reaches the end of a record without finding values for all fields, variables without values are set to missing.

Viewing the Output

Partial SAS Log

NOTE: 5 records were read from the infile "S:\workshop\phone.csv".

The minimum record length was 31.

The maximum record length was 44.

NOTE: The data set WORK.CONTACTS has 5 observations and 3 variables.

PROC PRINT Output

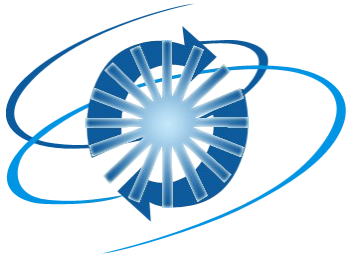
Name	Phone	Mobile
James Kvarniq	(704) 293-8126	(701) 281-8923
Sandrina Stephano	(919) 871-7830	
Cornelia Krah1	(212) 891-3241	(212) 233-5413
Karen Ballinger	(714) 344-4321	
Elke Wallstab	(910) 763-5561	(910) 545-3421

INFILE Options

```
INFILE "raw-data-file" <DLM=> <DSD> <MISSOVER>;
```

Option	Description
DLM=	Specifies an alternate delimiter.
DSD	Sets the default delimiter to a comma, treats consecutive delimiters as missing values, and allows embedded delimiters when the data value is enclosed in quotation marks.
MISSOVER	Sets variables to missing if the end of the record is reached before finding values for all fields.





Exercise

This exercise reinforces the concepts discussed previously.



Chapter Review

The graphic features the text 'Chapter Review' in a large, bold, blue sans-serif font. To the right of the text is a vertical stack of three circles. The middle circle contains a blue checkmark, while the top and bottom circles are empty. The entire text and icon set is enclosed within a large, stylized blue swoosh that curves around the top and bottom, with a slight gap on the right side.



1. In the first iteration of this program, SAS does the following:
 - reads a record from the raw data file into the input buffer
 - scans the input buffer and copies the values to the PDV
 - writes the values to the output data set
 - reinitializes the input buffer
 - reads the next record from the raw data file

```
data work.profit;  
    infile 'c:\mydata\income.csv' dlm=',';  
    input Amount SalesRep $ Customer $;  
run;
```

- True
- False

1. In the first iteration of this program, SAS does the following:
 - reads a record from the raw data file into the input buffer
 - scans the input buffer and copies the values to the PDV
 - writes the values to the output data set
 - reinitializes the input buffer
 - reads the next record from the raw data file.

```
data work.profit;  
    infile 'c:\mydata\income.csv' dlm=',';  
    input Amount SalesRep $ Customer $;  
run;
```

— True

☒ False

2. Which INFILE statement correctly specifies the raw data file shown here?

Partial **salestotals.dat**

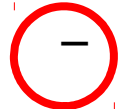
	1	1	2	2	3	3	4	4	5
1---	5---	0---	5---	0---	5---	0---	5---	0---	5---
14528*	instore*	06/15/2008*	215.65*	1650072*	red				
14529*	online*	06/15/2008*	183.98*	1650039*	white				
14530*	online*	06/16/2008*	107.50*	1650450*	green				
14531*	instore*	06/17/2008*	350.78*	1652903*	graphite				

- infile 'c:\mydata\salestotals.dat';
- infile 'c:\mydata\salestotals.dat' dlm=*;
- infile 'c:\mydata\salestotals.dat' dlm=',';
- infile 'c:\mydata\salestotals.dat' dlm='*';

2. Which INFILE statement correctly specifies the raw data file shown here?

Partial **salestotals.dat**

	1	1	2	2	3	3	4	4	5	
1---	5---	0---	5---	0---	5---	0---	5---	0---	5---	0
14528*instore*06/15/2008*215.65*1650072*red										
14529*online*06/15/2008*183.98*1650039*white										
14530*online*06/16/2008*107.50*1650450*green										
14531*instore*06/17/2008*350.78*1652903*graphite										

- infile 'c:\mydata\salestotals.dat';
- infile 'c:\mydata\salestotals.dat' dlm=*;
- infile 'c:\mydata\salestotals.dat' dlm=',';
-  - infile 'c:\mydata\salestotals.dat' dlm='*';

3. Which of the following INPUT statements creates the data set shown here, assuming that the DATA step does not contain a LENGTH statement?

Partial SAS Data Set **customers**

Customer_ID	Last_Name	First_Name	Total_Sales
123049	Kim	Jason	545
123050	Weston	Ingrid	832

- input Customer_ID \$ Last_Name \$ First_Name \$ Total_Sales;
- input customer_id \$ last_name \$ first_name \$ total_sales;
- input Last_Name \$ First_Name \$ Total_Sales Customer_ID \$;

3. Which of the following INPUT statements creates the data set shown here, assuming that the DATA step does not contain a LENGTH statement?

Partial SAS Data Set **customers**

Customer_ID	Last_Name	First_Name	Total_Sales
123049	Kim	Jason	545
123050	Weston	Ingrid	832

- input Customer_ID \$ Last_Name \$ First_Name \$ Total_Sales;
- input customer_id \$ last_name \$ first_name \$ total_sales;
- input Last_Name \$ First_Name \$ Total_Sales Customer_ID \$;

4. The INPUT statement below correctly reads this space-delimited raw data file.

	1	1	2	2	3	3
1	---	5	---	0	---	5
Donny	5MAY2008	25	FL	\$43,123.50		
Margaret	20FEB2008	43	NC	65,150		

```
input name $ hired date9. age state $  
       salary comma10.;
```

- True
- False

4. The INPUT statement below correctly reads this space-delimited raw data file.

	1	1	2	2	3	3
1	---	5	---	0	---	5
Donny	5MAY2008	25	FL	\$43,123.50		
Margaret	20FEB2008	43	NC	65,150		

```
input name $ hired date9. age state $  
       salary comma10.;
```

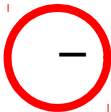
— True

☒ False

5. By default, SAS creates character variables with a length of _____ bytes for list input.

- 6
- 8
- 10
- 12

5. By default, SAS creates character variables with a length of _____ bytes for list input.

- 6
- - 8
- 10
- 12

6. Which of the following values can SAS store in a character variable that has a length of 8 bytes?
- Sales Manager
 - Regional Manager
 - 12036578
 - \$123,293.50
 - 06/15/2008

6. Which of the following values can SAS store in a character variable that has a length of 8 bytes?
- Sales Manager
 - Regional Manager
 - ☒ - 12036578
 - \$123,293.50
 - 06/15/2008

7. To explicitly define the length of a variable read from a raw data file, you use a LENGTH statement after the INPUT statement in a DATA step.

- | True
- | False

7. To explicitly define the length of a variable read from a raw data file, you use a LENGTH statement after the INPUT statement in a DATA step.

—| True

☒ False

8. Which INFILE statement correctly specifies the raw data file shown here?

Partial **sales.dat**

	1	1	2	2	3	3	4	4	5			
1	---	5	---	0	---	5	---	0	---	5	---	0
14528	*	instore	*	06/15/2008	*	215.65	*	red				
14529	*	online	*	06/15/2008	*	183.98	*	1650039	*	white		
14530	*	*	06/16/2008	*	107.50	*	1650450	*	green			
14531	*	instore	*	06/17/2008	*	350.78	*	1652903	*	graphite		

- infile 'c:\mydata\sales.dat';
- infile 'c:\mydata\sales.dat' dsd dlm='*';
- infile 'c:\mydata\sales.dat' dlm=*;
- infile 'c:\mydata\sales.dat' dlm='*';

8. Which INFILE statement correctly specifies the raw data file shown here?

Partial **sales.dat**

	1	1	2	2	3	3	4	4	5			
1	---	5	---	0	---	5	---	0	---	5	---	0
14528	*	instore	*	06/15/2008	*	215.65	*	red				
14529	*	online	*	06/15/2008	*	183.98	*	1650039	*	white		
14530	*		*	06/16/2008	*	107.50	*	1650450	*	green		
14531	*	instore	*	06/17/2008	*	350.78	*	1652903	*	graphite		

- infile 'c:\mydata\sales.dat';
- ☒ - infile 'c:\mydata\sales.dat' dsd dlm='*';
- infile 'c:\mydata\sales.dat' dlm='*';
- infile 'c:\mydata\sales.dat' dlm='*';

9. Which of the following statements specifies in-stream data, or the lines of data that you enter directly in a DATA step?
- DATALINES
 - INFILE
 - INPUT
 - INSTREAM


9. Which of the following statements specifies in-stream data, or the lines of data that you enter directly in a DATA step?

- ☒ – DATALINES
- INFILE
- INPUT
- INSTREAM

10. Which of the following statements cannot be used in a DATA step that reads a raw data file as input?

- KEEP
- IF
- FORMAT
- WHERE

10. Which of the following statements cannot be used in a DATA step that reads a raw data file as input?

- KEEP
- IF
- FORMAT
-  – WHERE