

Cleaning Data with Text Functions

The data that you obtain from different sources may not be in a form ready for analysis. In this chapter, you will understand how to prepare your data that is in the form of text for analysis.

Initially, you need to clean the data. Data cleaning includes removing unwanted characters from text. Next, you need to structure the data in the form you require for further analysis. You can do the same by –

- Finding required text patterns with the text functions.
- Extracting data values from text.
- Formatting data with text functions.
- Executing data operations with the text functions.

Removing Unwanted Characters from Text

When you import data from another application, it can have nonprintable characters and/or excess spaces. The excess spaces can be –

- leading spaces, and/or
- extra spaces between words.

If you sort or analyze such data, you will get erroneous results.

Consider the following example –

Product Data	
54482100AFES	CONTROLLER SERVER 1TB H 304.00
54482100ICP9	DESKTOP UNIT 225.00
54482700BAAS	DESKTOP WINDOWS 8.1 SERVER 2302.00
54482600BAAS	DESKTOP WINDOWS 8.1 WKST 355.00
54482100BAAS	DESKTOP WINDOWS 10 182.00
54482200BAAS	DESKTOP WINDOWS DESKTOP OS 255.00
54482500BAAS	DESKTOP WINDOWS OS 354.00
54483000BAAS	MINITOWER NO OS 1840.00
54483000KEBB	MINI TOWER 2550.00

This is the raw data that you have obtained on product information containing the Product ID, Product description and the price. The character "|" separates the field in each row.

When you import this data into Excel worksheet, it looks as follows –

A	B	C
	Product Data	
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	
4	54482100ICP9 DESKTOP UNIT 225.00	
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	
7	54482100BAAS DESKTOP WINDOWS 10 182.00	
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	
9	54482500BAAS DESKTOP WINDOWS OS 354.00	
10	54483000BAAS MINITOWER NO OS 1840.00	
11	54483000KEBB MINI TOWER 2550.00	

As you observe, the entire data is in a single column. You need to structure this data to perform data analysis. However, initially you need to clean the data.

You need to remove any nonprintable characters and excess spaces that might be present in the data. You can use the CLEAN function and TRIM function for this purpose.

S.No.	Function & Description
1.	CLEAN Removes all nonprintable characters from text
2.	TRIM Removes spaces from text

- Select the Cells C3 – C11.
- Type =TRIM (CLEAN (B3)) and then press CTRL + Enter.

The formula is filled in the cells C3 – C11.

A	B	C
	Product Data	
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	=TRIM(CLEAN(B3))
4	54482100ICP9 DESKTOP UNIT 225.00	=TRIM(CLEAN(B4))
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	=TRIM(CLEAN(B5))
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	=TRIM(CLEAN(B6))
7	54482100BAAS DESKTOP WINDOWS 10 182.00	=TRIM(CLEAN(B7))
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	=TRIM(CLEAN(B8))
9	54482500BAAS DESKTOP WINDOWS OS 354.00	=TRIM(CLEAN(B9))
10	54483000BAAS MINITOWER NO OS 1840.00	=TRIM(CLEAN(B10))
11	54483000KEBB MINI TOWER 2550.00	=TRIM(CLEAN(B11))

The result will be as shown below –

A	B	C
	Raw Data	Nonprintable Characters and Excess Spaces removed
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	54482100AFES CONTROLLER SERVER 1TB H 304.00
4	54482100ICP9 DESKTOP UNIT 225.00	54482100ICP9 DESKTOP UNIT 225.00
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00
7	54482100BAAS DESKTOP WINDOWS 10 182.00	54482100BAAS DESKTOP WINDOWS 10 182.00
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00
9	54482500BAAS DESKTOP WINDOWS OS 354.00	54482500BAAS DESKTOP WINDOWS OS 354.00
10	54483000BAAS MINITOWER NO OS 1840.00	54483000BAAS MINITOWER NO OS 1840.00
11	54483000KEBB MINI TOWER 2550.00	54483000KEBB MINI TOWER 2550.00

Finding required Text Patterns with the Text Functions

To structure your data, you might have to do certain Text Pattern matching based on which you can extract the Data Values. Some of the Text Functions that are useful for this purpose are –

S.No.	Function & Description
1.	EXACT Checks to see if two text values are identical
2.	FIND Finds one text value within another (case-sensitive)
3.	SEARCH Finds one text value within another (not case-sensitive)

Extracting Data Values from Text

You need to extract the required data from text in order to structure the same. In the above example, say, you need to place the data in three columns – ProductID, Product_Description and Price.

You can extract data in one of the following ways –

- Extracting Data Values with Convert Text to Columns Wizard
- Extracting Data Values with Text Functions
- Extracting Data Values with Flash Fill

Extracting Data Values with Convert Text to Columns Wizard

You can use the **Convert Text to Columns Wizard** to extract Data Values into Excel columns if your fields are –

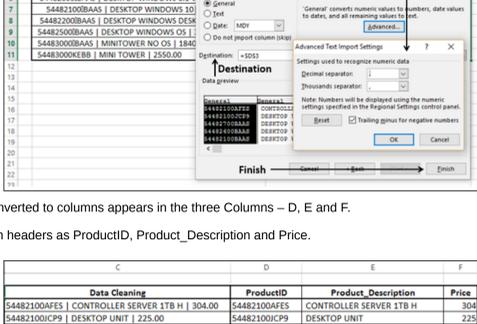
- Delimited by a character, or
- Aligned in columns with spaces between each field.

In the above example, the fields are delimited by the character "|". Hence, you can use the **Convert Text to Columns wizard**.

- Select the data.
- Click on **Text to Columns** in the **Data Tools** group under **Data Tab** on the Ribbon.

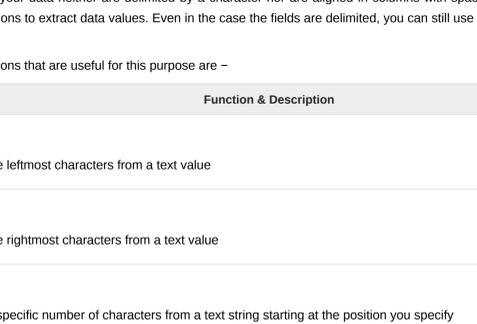
Step 1 – Convert Text to Columns Wizard – Step 1 of 3 appears.

- Select **Delimited**.
- Click **Next**.



Step 2 – Convert Text to Columns Wizard – Step 2 of 3 appears.

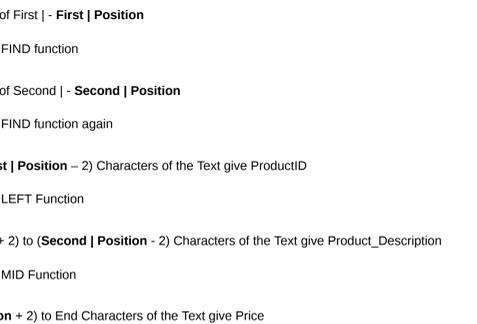
- Under **Delimiters**, select **Other**.
- In the box next to **Other**, type the character |
- Click **Next**.



Step 3 – Convert Text to Columns Wizard – Step 3 of 3 appears.

In this screen, you can select each column of your data in the wizard and set the format for that column.

- For **Destination**, select the cell D3.
- You can click **Advanced**, and set **Decimal Separator** and **Thousands Separator** in the **Advanced Text Import Settings** dialog box that appears.
- Click **Finish**.



Your data, which is converted to columns appears in the three Columns – D, E and F.

- Name the Column headers as ProductID, Product_Description and Price.

	C	D	E	F	
	Data Cleaning		ProductID	Product_Description	Price
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	54482100AFES	CONTROLLER SERVER 1TB H	304	
4	54482100ICP9 DESKTOP UNIT 225.00	54482100ICP9	DESKTOP UNIT	225	
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	54482700BAAS	DESKTOP WINDOWS 8.1 SERVER	2302	
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	54482600BAAS	DESKTOP WINDOWS 8.1 WKST	355	
7	54482100BAAS DESKTOP WINDOWS 10 182.00	54482100BAAS	DESKTOP WINDOWS 10	182	
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255	54482200BAAS	DESKTOP WINDOWS DESKTOP OS	255	
9	54482500BAAS DESKTOP WINDOWS OS 354.00	54482500BAAS	DESKTOP WINDOWS OS	354	
10	54483000BAAS MINITOWER NO OS 1840.00	54483000BAAS	MINITOWER NO OS	1840	
11	54483000KEBB MINI TOWER 2550.00	54483000KEBB	MINI TOWER	2550	

Extracting Data Values with Text Functions

Suppose the fields in your data neither are delimited by a character nor are aligned in columns with spaces between each field, you can use text functions to extract data values. Even in the case the fields are delimited, you can still use text functions to extract data.

Some of the text functions that are useful for this purpose are –

S.No.	Function & Description
1.	LEFT Returns the leftmost characters from a text value
2.	RIGHT Returns the rightmost characters from a text value
3.	MID Returns a specific number of characters from a text string starting at the position you specify
4.	LEN Returns the number of characters in a text string

You can also combine two or more of these text functions as per the data you have at hand, to extract the required data values. For example, using a combination of LEFT, RIGHT and VALUE functions or using a combination of FIND, LEFT, LEN and MID functions.

In the above example,

- All the characters left to the first | give the name ProductID.
- All the characters right to the second | give the name Price.
- All the characters that lie between the first | and second | give the name Product_Description.
- Each | has a space before and after.

Observing this information, you can extract the data values with the following steps –

- Find the Position of First | - **First | Position**
 - You can use FIND function
- Find the Position of Second | - **Second | Position**
 - You can use FIND function again
- Beginning to (First | Position – 2) Characters of the Text give ProductID
 - You can use LEFT Function
- (First | Position + 2) to (Second | Position - 2) Characters of the Text give Product_Description
 - You can use MID Function
- (Second | Position + 2) to End Characters of the Text give Price
 - You can use RIGHT Function

	B	C	D	E	F	G	
	Product Data		First Position	Second Position	ProductID	Product_Description	Price
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	=FIND(" ",B3)	=FIND(" ",B3)	=LEFT(B3,C3-D3)	=MID(B3,E3-F3)	=RIGHT(B3,G3)	
4	54482100ICP9 DESKTOP UNIT 225.00	=FIND(" ",B4)	=FIND(" ",B4)	=LEFT(B4,C4-D4)	=MID(B4,E4-F4)	=RIGHT(B4,G4)	
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	=FIND(" ",B5)	=FIND(" ",B5)	=LEFT(B5,C5-D5)	=MID(B5,E5-F5)	=RIGHT(B5,G5)	
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	=FIND(" ",B6)	=FIND(" ",B6)	=LEFT(B6,C6-D6)	=MID(B6,E6-F6)	=RIGHT(B6,G6)	
7	54482100BAAS DESKTOP WINDOWS 10 182.00	=FIND(" ",B7)	=FIND(" ",B7)	=LEFT(B7,C7-D7)	=MID(B7,E7-F7)	=RIGHT(B7,G7)	
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	=FIND(" ",B8)	=FIND(" ",B8)	=LEFT(B8,C8-D8)	=MID(B8,E8-F8)	=RIGHT(B8,G8)	
9	54482500BAAS DESKTOP WINDOWS OS 354.00	=FIND(" ",B9)	=FIND(" ",B9)	=LEFT(B9,C9-D9)	=MID(B9,E9-F9)	=RIGHT(B9,G9)	
10	54483000BAAS MINITOWER NO OS 1840.00	=FIND(" ",B10)	=FIND(" ",B10)	=LEFT(B10,C10-D10)	=MID(B10,E10-F10)	=RIGHT(B10,G10)	
11	54483000KEBB MINI TOWER 2550.00	=FIND(" ",B11)	=FIND(" ",B11)	=LEFT(B11,C11-D11)	=MID(B11,E11-F11)	=RIGHT(B11,G11)	

The result will be as shown below –

	B	C	D	E	F	G	
	Product Data		First Position	Second Position	ProductID	Product_Description	Price
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	14	40	54482100AFES	CONTROLLER SERVER 1TB H	304.00	
4	54482100ICP9 DESKTOP UNIT 225.00	14	29	54482100ICP9	DESKTOP UNIT	225.00	
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	14	45	54482700BAAS	DESKTOP WINDOWS 8.1 SERVER	2302.00	
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	14	41	54482600BAAS	DESKTOP WINDOWS 8.1 WKST	355.00	
7	54482100BAAS DESKTOP WINDOWS 10 182.00	14	34	54482100BAAS	DESKTOP WINDOWS 10	182.00	
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	14	40	54482200BAAS	DESKTOP WINDOWS DESKTOP OS	255.00	
9	54482500BAAS DESKTOP WINDOWS OS 354.00	14	35	54482500BAAS	DESKTOP WINDOWS OS	354.00	
10	54483000BAAS MINITOWER NO OS 1840.00	14	31	54483000BAAS	MINITOWER NO OS	1840.00	
11	54483000KEBB MINI TOWER 2550.00	14	27	54483000KEBB	MINITOWER	2550.00	

You can observe that the values in the price column are text values. To perform calculations on these values, you have to format the corresponding cells. You can look at the section given below to understand formatting text.

Extracting Data Values with Flash Fill

Using Excel **Flash Fill** is another way to extract data values from text. However, this works only when Excel is able to find a pattern in the data.

Step 1 – Create three columns for ProductID, Product_Description and Price next to the data.

A	B	C	D	E
	Product Data		ProductID	Product_Description
3	54482100AFES CONTROLLER SERVER 1TB H 304.00			
4	54482100ICP9 DESKTOP UNIT 225.00			
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00			
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00			
7	54482100BAAS DESKTOP WINDOWS 10 182.00			
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00			
9	54482500BAAS DESKTOP WINDOWS OS 354.00			
10	54483000BAAS MINITOWER NO OS 1840.00			
11	54483000KEBB MINI TOWER 2550.00			

Step 2 – Copy and paste the values for C3, D3 and E3 from B3.

A	B	C	D	E
	Product Data		ProductID	Product_Description
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	54482100AFES	CONTROLLER SERVER 1TB H	304
4	54482100ICP9 DESKTOP UNIT 225.00	54482100ICP9	DESKTOP UNIT	225
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	54482700BAAS	DESKTOP WINDOWS 8.1 SERVER	2302
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	54482600BAAS	DESKTOP WINDOWS 8.1 WKST	355
7	54482100BAAS DESKTOP WINDOWS 10 182.00	54482100BAAS	DESKTOP WINDOWS 10	182
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	54482200BAAS	DESKTOP WINDOWS DESKTOP OS	255
9	54482500BAAS DESKTOP WINDOWS OS 354.00	54482500BAAS	DESKTOP WINDOWS OS	354
10	54483000BAAS MINITOWER NO OS 1840.00	54483000BAAS	MINITOWER NO OS	1840
11	54483000KEBB MINI TOWER 2550.00	54483000KEBB	MINI TOWER	2550

Step 3 – Select cell C3 and click Flash Fill in the Data Tools group on the Data tab. All the values for ProductID get filled.

A	B	C	D	E	F	G
	Product Data		ProductID	Product_Description	Price	
3	54482100AFES CONTROLLER SERVER 1TB H 304.00	54482100AFES	CONTROLLER SERVER 1TB H	304		
4	54482100ICP9 DESKTOP UNIT 225.00	54482100ICP9	DESKTOP UNIT	225		
5	54482700BAAS DESKTOP WINDOWS 8.1 SERVER 2302.00	54482700BAAS	DESKTOP WINDOWS 8.1 SERVER	2302		
6	54482600BAAS DESKTOP WINDOWS 8.1 WKST 355.00	54482600BAAS	DESKTOP WINDOWS 8.1 WKST	355		
7	54482100BAAS DESKTOP WINDOWS 10 182.00	54482100BAAS	DESKTOP WINDOWS 10	182		
8	54482200BAAS DESKTOP WINDOWS DESKTOP OS 255.00	54482200BAAS	DESKTOP WINDOWS DESKTOP OS	255		
9	54482500BAAS DESKTOP WINDOWS OS 354.00	54482500BAAS	DESKTOP WINDOWS OS	354		
10	54483000BAAS MINITOWER NO OS 1840.00	54483000BAAS				