

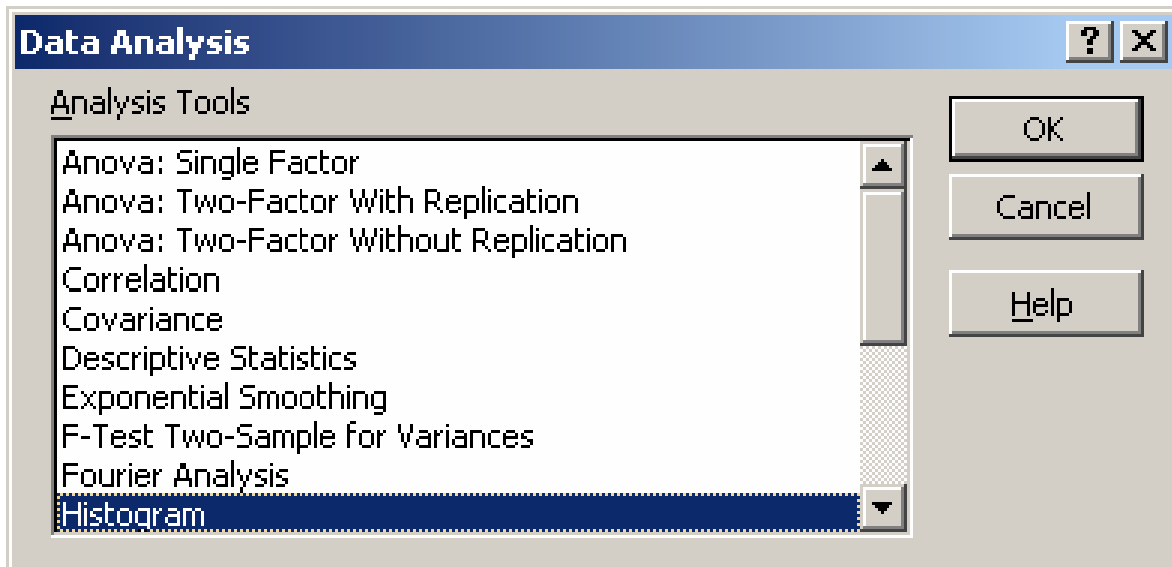
Using Excel's Histogram command

This handout briefly discusses the *Histogram* command in the Data Analysis Add-In of Excel. It assumes that the Data Analysis Add-In has been installed. We will use the Housing Data as the example data set. Our goal is to generate a histogram for the *square footage* of all houses in the data set. From the Command Menu, choose *Tools*, then *Data Analysis* as in the graphic below:

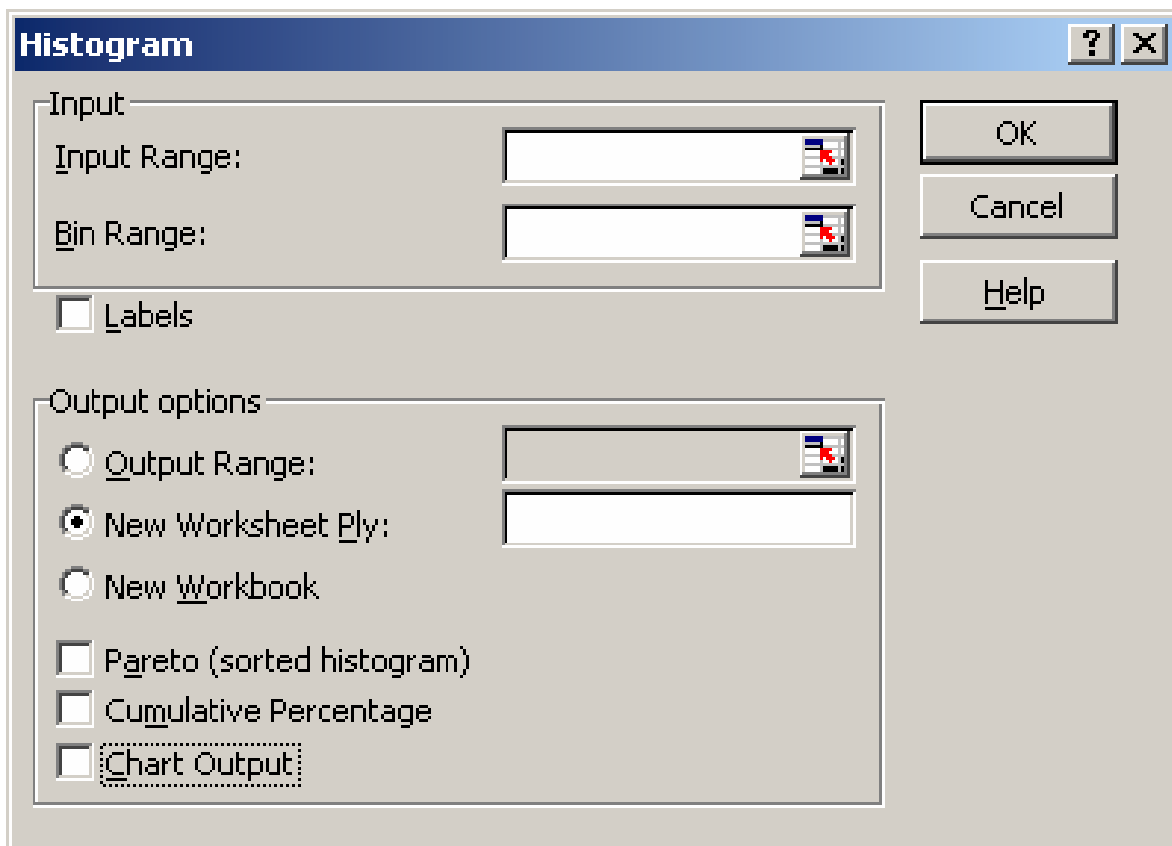
The screenshot shows the Microsoft Excel interface with the 'Tools' menu open and 'Data Analysis...' selected. The spreadsheet contains housing data with columns for house number, square footage, bedrooms, bathrooms, price, and school district. A text box on the right explains the variables:

- SOFT** is the variable measuring the total square feet in the house.
- BEDS** and **BATHS** are number of bedrooms and bathrooms, respectively.
- HEAT** and **STYLE** are categorical variables. **HEAT** takes on the value of zero for oil heating and one for electric heat. **STYLE** is the architectural style of the home: zero indicates a Cape Cod style, one indicates a two-story house and two indicates that the house is a ranch-style house.
- GARAGE** is the number of cars that can fit into the garage. (We don't know if it is an attached garage.)
- AGE** is the age of the home in years.
- FIRE** and **BASEMENT** indicate the presence (one) or absence (zero) of a fireplace or basement.
- PRICE** is the selling price of the house in thousands of dollars.
- SCHOOL** is the school district (0=Plum Ridge school district; 1=Apple Valley school district). Apple Valley seems to be viewed as the preferred school district.

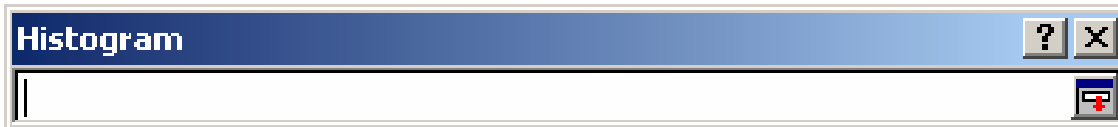
The next dialog box for *Data Analysis* appears, you choose the option *Histogram* from the list and press OK.



The following dialog box then appears:



Input Range is the range in the Excel spreadsheet that contains the numerical values for which we want to generate a histogram. By pressing the *red arrow button* on the right, Excel responds with:



and we can now freely select the range in the spreadsheet directly. You may also type the starting and ending address of the range manually (e.g.: B1:B109). Pressing *Enter* on the keyboard will take you to the *Bin Range* field. *Bins* are classes in Excel language. If you leave this field blank, Excel will automatically compute what it thinks are the “best” classes. But, it will likely end up with class boundaries that you won’t like, in this case, e.g., the first class starts at 816 sq.ft. and ends at 1015.3 sq. ft., whereas we would rather like to see class boundaries like 800 through 1000, 1000 through 1200, etc. In this case we need to specify classes for the bins somewhere in the spreadsheet, something like this (the bins could be in another sheet as well, different from the sheet where the source data resides):

Microsoft Excel - HouseData

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	house	sqft	beds	baths	heat	style	garage	basement	age	fire	price	school							
2	1	1,238	3	2	0	0	1	1	12	1	69,900	1							
3	2	1,707	3	2	1	0	2	0	13	1	64,000	0							
4	3	1,296	4	2	0	0	2	1	17	0	66,500	0							
5	4	1,320	3	2	0	0	2	1	11	1	66,500	0							
6	5	1,210	3	2	0	0	1	0	6	1	66,900	0							
7	6	1,296	3	2	0	0	2	1	17	1	68,000	0							
8	7	1,765	3	2	0	0	2	1	20	0	68,500	0							
9	8	1,725	4	3	0	0	2	1	12	0	69,000	0							
10	9	1,794	4	2	0	0	2	1	18	0	70,950	0							
11	10	1,294	3	2	0	0	2	0	13	1	71,000	0							
12	11	1,372	3	2	0	0	2	1	9	0	72,692	1							
13	12	1,162	3	2	0	0	1	0	8	1	72,601	0							
14	13	1,996	4	2	0	0	2	0	13	1	75,207	1							
15	14	1,764	4	2	1	0	2	1	13	1	76,000	0							
16	15	1,416	3	2	0	0	2	0	8	0	76,000	1							
17	16	1,730	4	2	0	0	2	0	15	1	77,500	0							
18	17	1,392	3	2	0	0	2	1	8	1	79,900	1							
19	18	1,664	3	3	0	0	2	0	11	1	79,900	0							
20	19	1,332	3	2	0	0	2	1	14	0	81,000	1							
21	20	1,752	3	3	0	0	2	0	18	1	82,800	0							
22	21	2,167	3	3	1	0	2	1	13	1	84,900	0							
23	22	1,664	3	2	0	0	2	0	9	1	85,000	0							
24	23	1,973	4	3	0	0	2	0	13	1	86,000	0							
25	24	1,384	3	2	0	0	2	0	5	1	89,280	0							
26	25	1,431	3	2	0	0	2	1	7	1	89,900	1							
27	26	1,960	5	3	0	0	2	0	13	1	90,000	0							
28	27	1,452	3	2	0	0	2	1	4	1	92,000	0							
29	28	1,829	3	2	0	0	2	0	10	1	92,439	1							
30	29	1,652	4	3	0	0	2	1	7	1	94,646	1							
31	30	1,516	3	2	0	0	2	1	10	1	97,293	1							
32	31	1,998	4	3	0	0	2	1	17	1	98,100	1							
33	32	1,984	4	3	0	0	2	1	9	1	98,149	0							
34	33	1,840	3	3	0	0	2	1	9	1	105,000	0							
35	34	1,823	4	3	0	0	2	1	3	1	110,000	1							
36	35	2,150	5	3	0	0	2	1	12	1	111,900	0							
37	36	2,096	3	3	0	0	2	1	9	1	113,000	1							

* **SQFT** is the variable measuring the total square feet in the house.
 * **BEDS** and **BATHS** are number of bedrooms and bathrooms, respectively.
 * **HEAT** and **STYLE** are categorical variables.
HEAT takes on the value of zero for oil heating and one for electric heat.
STYLE is the architectural style of the home: zero indicates a Cape Cod sty indicates a two-story house and two indicates that the house is a ranch-sty
 * **GARAGE** is the number of cars that can fit into the garage. (We don't know is an attached garage.)
 * **AGE** is the age of the home in years.
 * **FIRE** and **BASEMENT** indicate the presence (one) or absence (zero) of at le fireplace or basement.
 * **PRICE** is the selling price of the house in thousands of dollars.
 * **SCHOOL** is the school district (0=Plum Ridge school district; 1=Apple Vall district). Apple Valley seems to be viewed as the preferred school district.

Bins

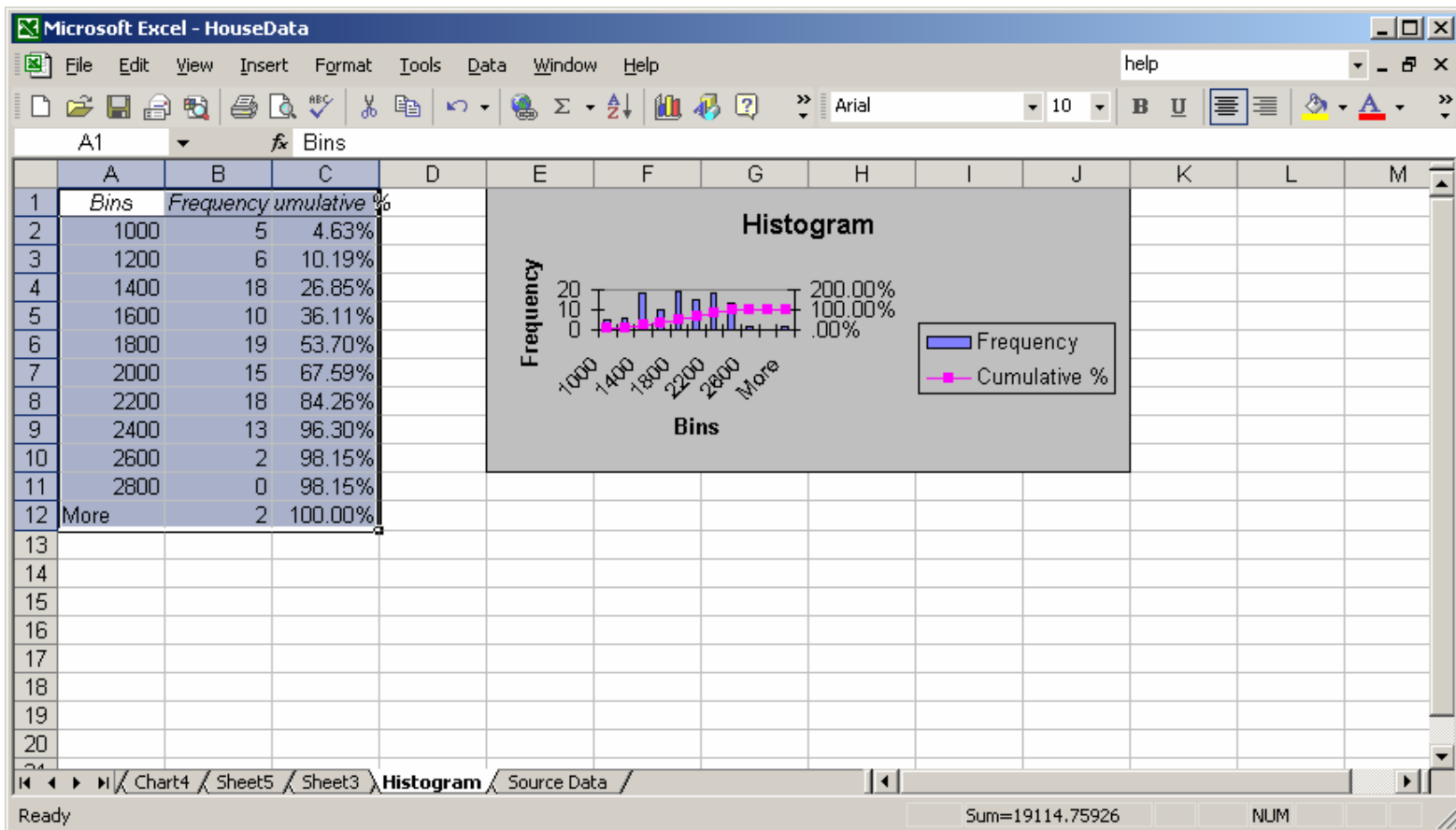
1000
1200
1400
1600
1800
2000
2200
2400
2600
2800

Excel will automatically create a last *bin*, and call it “More” for the values that are higher than 2800. By pressing the red arrow button on the dialog box at the *Bin Range* field,

we can select the cells that contain our classes (or again, type in range manually). Be sure to select the cell above the first one as well in your range if you check the *Labels* box, otherwise Excel thinks that your first value ("1000") is actually a label and not the ending boundary of your first class.

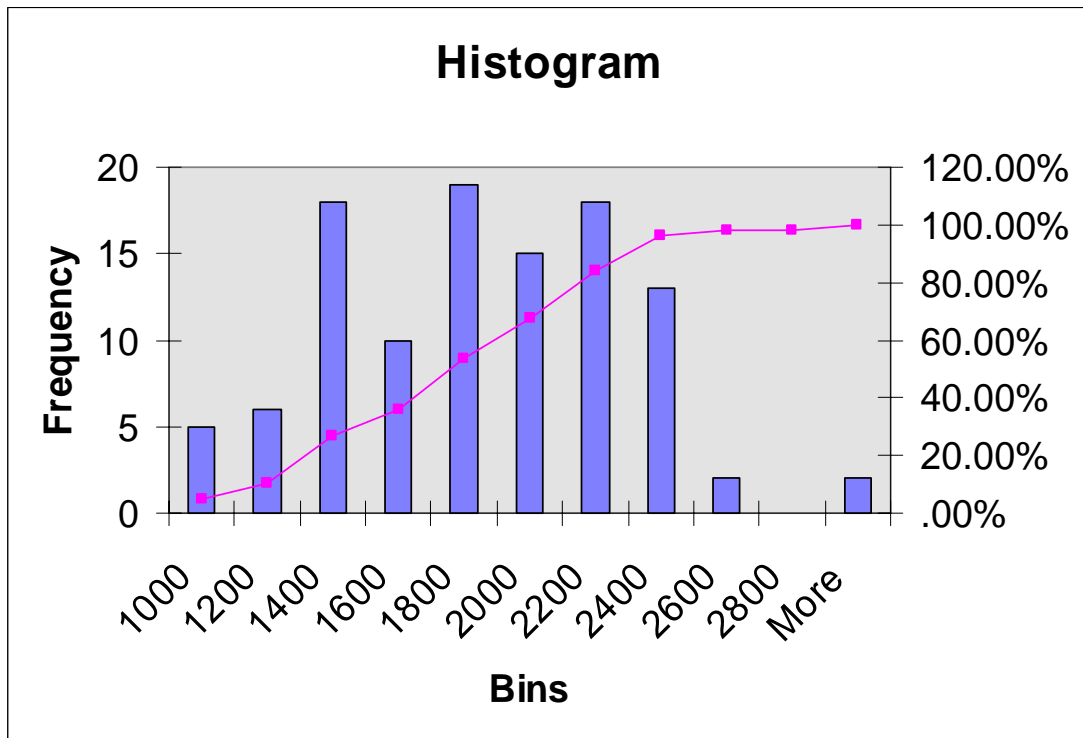
For the output options, we can decide to put the histogram on the same spreadsheet (first option), in which case we need to specify where the results should go (the field "*Output Range*"), or put the results in a new sheet in the same book (we may give it a name such as "histogram" in our example) or we can choose to put the results in a new workbook.

The last two checkboxes allow you to compute a *cumulative frequency distribution* and have the results displayed in a *chart* (graph), which was our goal from the outset. After pressing the *OK* button, the new sheet called "histogram" is created which looks as follows:



The column *Frequency* (and *Cumulative %*) needs to be interpreted as the number of observations *less than or equal to* the value under *Bins*, i.e., there are 5 observations with a square footage less than or equal to 1000, 6 observations with values between 1000 and 1200, etc.

The charted histogram is probably a little too small, but by clicking on it, you can drag it and stretch it and if you want to do some more basic formatting, you may get it to look like this:



Wow, that's one good looking histogram with very little effort!