

MATLAB Reference 3: Descriptive Statistics & Histograms

Descriptive Statistics

| Command & syntax | Purpose |
|------------------------|--|
| <code>min(A)</code> | If A is an array, returns the minimum value in the array A. If A is a matrix, returns the minimum value from each column. |
| <code>max(A)</code> | If A is an array, returns the maximum value in the array A. If A is a matrix, returns the minimum value from each column. |
| <code>range(A)</code> | If A is an array, returns the difference between the minimum and maximum values in the array A. If A is a matrix, returns the differences between the minimum and maximum values in each column. |
| <code>length(A)</code> | Returns the number of elements in the longest dimension of A. |
| <code>size(A)</code> | Returns the number of elements in each dimension of A. |

Measures of central tendency of the data, or what is typical of the distribution of the data

| Command & syntax | Purpose |
|------------------------|---|
| <code>mean(A)</code> | If A is an array, returns the average or arithmetic mean . If A is a matrix, returns the average or arithmetic mean of each column. |
| <code>median(A)</code> | If A is an array, returns the median value . If A is a matrix, returns the median of each column. |
| <code>mode(A)</code> | If A is an array, most frequently occurring value . If A is a matrix, returns the most frequently occurring value of each column. |

Measures of spread of the data

| Command & syntax | Purpose |
|---------------------|---|
| <code>var(A)</code> | Returns the sample variance of all the values in A Note: The variance computed by MATLAB's <code>var(A)</code> is different from the variance of Python's <code>numpy</code> or <code>scipy</code> modules. Read the note below for further information. |

| | |
|---------------------|--|
| <code>std(A)</code> | Returns the sample standard deviation of all the values in A Note: The standard deviation computed by MATLAB's <code>std(A)</code> is different from the standard deviation of Python's <code>numpy</code> or <code>scipy</code> modules. |
|---------------------|--|

Difference of variance and standard deviation formulas in MATLAB, Python, Excel:

MATLAB's `var(A)` computes the sample variance: $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$, where n is the number of samples, x_i is the i -th data point and \bar{x} is the arithmetic mean of the sample.

Python's `numpy` and `scipy` modules compute by

```
import numpy; numpy.var(A) or
```

```
import scipy; scipy.var(A)
```

the population variance defined by $\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$, where N is the number of all data points in the population and x_i and \bar{x} are as before.

Note that the standard deviation is the (positive) square root of the variance. Use `std` instead of `var` in the above code to compute the respective standard deviation.

Why is there a difference? The difference comes from the Bessel correction [https://en.wikipedia.org/wiki/Bessel%27s_correction]. In short, if you use the population formula, your result is inaccurate if you look at only a sample. If you e.g. compute the standard deviation for 30 exam scores and there are only 30 students in the class, then you must use the population formulas. If you know there are 30 students in the class, but you only have the grades from 14, then you must use the sample formulas or your result is inaccurate.

How to compute the sample standard deviation in Python?

```
import numpy; numpy.std(A, ddof=1) # ddof means degrees of freedom
```

How to compute the population standard deviation in MATLAB?

```
std(A,1) % the 1 tells MATLAB to compute the population statistic
```

How to do it in Excel?

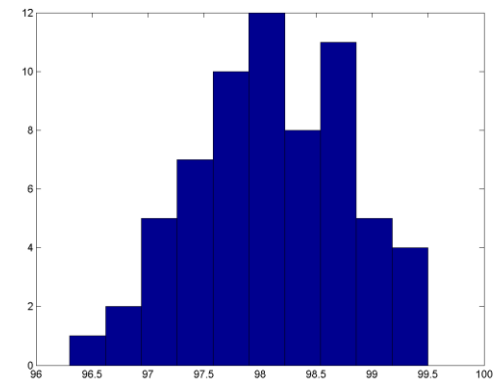
| | A | B | C | D | E | F | G | H |
|---|------------------------|---|---|---|---------------|-------------------------------|---|---|
| 1 | 1 | 2 | 3 | 4 | 5 | | | |
| 2 | =STDEV.S(A1:E1) | | | | 1.5811 | sample standard deviation | | |
| 3 | =STDEV.P(A1:E1) | | | | 1.4142 | population standard deviation | | |
| 4 | =VAR.S(A1:E1) | | | | 2.5 | sample variance | | |
| 5 | =VAR.P(A1:E1) | | | | 2 | population variance | | |

Histogram Commands

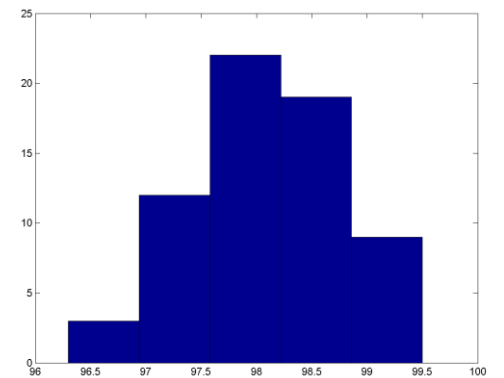
| Command & syntax | Purpose |
|--|--|
| <code>hist(temps)</code> | Returns a histogram of the dataset “temps” with 10 bins (default) |
| <code>hist(temps, X)</code> , where X is a scalar | Where X is a scalar defining the number of bins, returns a histogram of the dataset “temps” |
| <code>hist(temps, Y)</code> , where Y is a vector | Where Y is a vector that defines bin centers, returns a histogram of the dataset “temps” and distributes the data points among bins with centers specified in vector Y |
| <code>freq = hist(temps)</code> | Returns the number of data points in each bin for a default of 10 bins |
| <code>freq = hist(temps, X)</code> | Where X is the number of bins, returns the number of data points in each bin |
| <code>freq = hist(temps, Y)</code> | Where Y is a vector that defines bin centers, returns the number of data points in each bin |
| <code>[freq, centers] = hist(temps, X)</code> | Where X is the number of bins, returns the number of data points in each bin (freq) AND Returns the values of the bin centers (centers) |
| <code>centers_manual = [start: increment: end] hist(male_temps, centers_manual)</code> | Manually define bin centers |

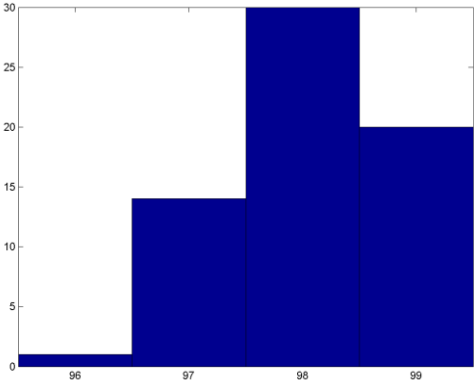
Example Histograms and Commands

```
hist(data)      % generate default histogram
```



```
nbr_bins=5      % specify number of bins  
hist(data,nbr_bins) % histogram with a specified number of bins
```



| | |
|---|--|
| <pre>centers =[96,97,98,99] % specify bin centers hist(data,centers) % histogram with a specified bin % centers</pre> |  <p>A histogram with four blue bars. The x-axis is labeled with bin centers 96, 97, 98, and 99. The y-axis represents frequency, ranging from 0 to 30. The bar heights are 3 for 96, 14 for 97, 30 for 98, and 20 for 99.</p> |
| <pre>freq=hist(data,nbr_bins) % determine # of data points in each % bin with number of bins specified</pre> | <pre>freq=[3 12 22 19 9]</pre> |
| <pre>freq=hist(data,centers) % determine # of data points in each % bin with bin centers specified</pre> | <pre>freq=[1 14 30 20]</pre> |
| <pre>[freq, centers] = hist(data, nbr_bins) % determine # of data % points in each bin and % the bin centers with % the number of bins % specified</pre> | <pre>centers= [96.62 97.26 97.9 98.54 99.18] freq=[3 12 22 19 9]</pre> |