

# AP Statistics – Lesson Notes - Chapter 6: Standard Deviation, Normal Model

## Comparing results from different datasets

Scores for college-bound students on the SAT and ACT tests are unimodal and symmetrical with the following means and standard deviations:

SAT:  $\bar{x} = 1500$ ,  $s = 250$       ACT:  $\bar{x} = 20.8$ ,  $s = 4.8$

Which of the following students has a better score?

Student 1: SAT score of 2030

Student 2: ACT score of 32

Both seem substantially above the mean. Standard deviation is a measure of spread, so what if we figured out how many standard deviations each is above its mean?

Student 1:

$$\frac{2030 - 1500}{250} = 2.12$$

std devs above mean

Student 2:

$$\frac{32 - 20.8}{4.8} = 2.33$$

std devs above mean

Determining the number of standard deviations a given data value,  $x$ , is above its mean is called calculating the data value's **z-score**:

$$z = \frac{x - \bar{x}}{s}$$

← how different this value is from the mean  
← average difference from the mean

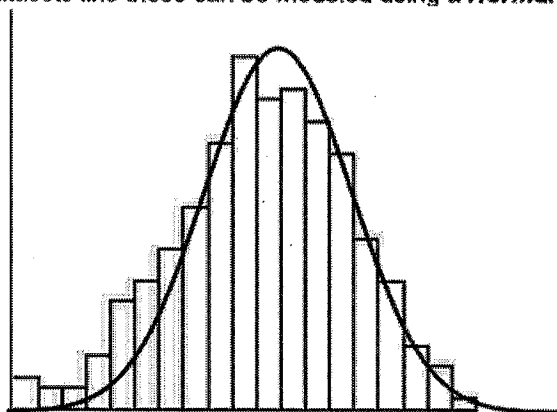
Because z-score has no units, it can be used to compare results from different datasets.

## Normal distribution model

How many standard deviations away from the mean does a data value have to be to be considered 'significantly different' than the mean?

In order to answer this, we need a way to model the distribution of the data values. Many datasets are unimodal and symmetrical, with most of the data grouped around the mean and lower frequency as you move away from the mean on both sides.

Datasets like these can be modeled using a **Normal Distribution Model**:



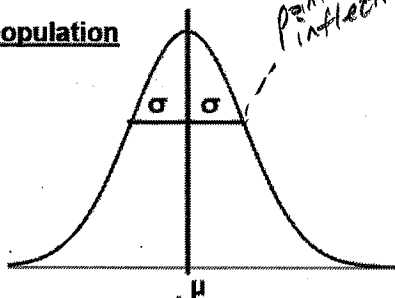
Don't need to know, but in case you're interested here is the probability density function model for a Normal distribution:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

(geogebra demo 'normalparameters.ggb')

**Normal distribution model**

The population



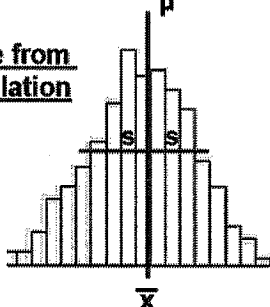
**Greek letters for populations**

$\mu$  = mean (mu 'mew')

$\sigma$  = standard (sigma) deviation

These are called the *parameters*.

A sample from the population



**English alphabet letters for samples**

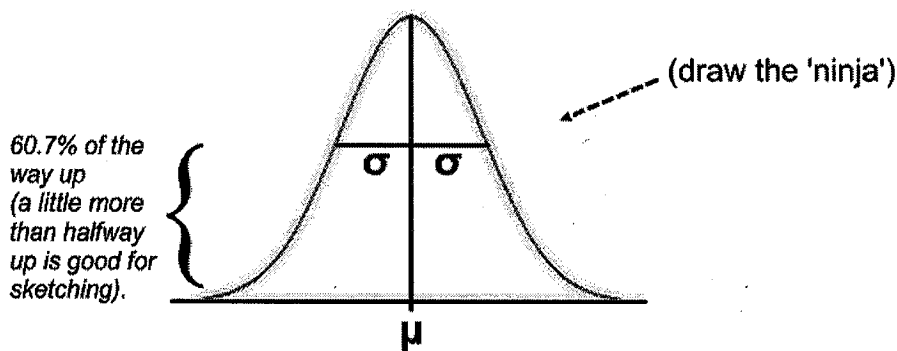
$\bar{x}$  = mean

$s$  = standard deviation

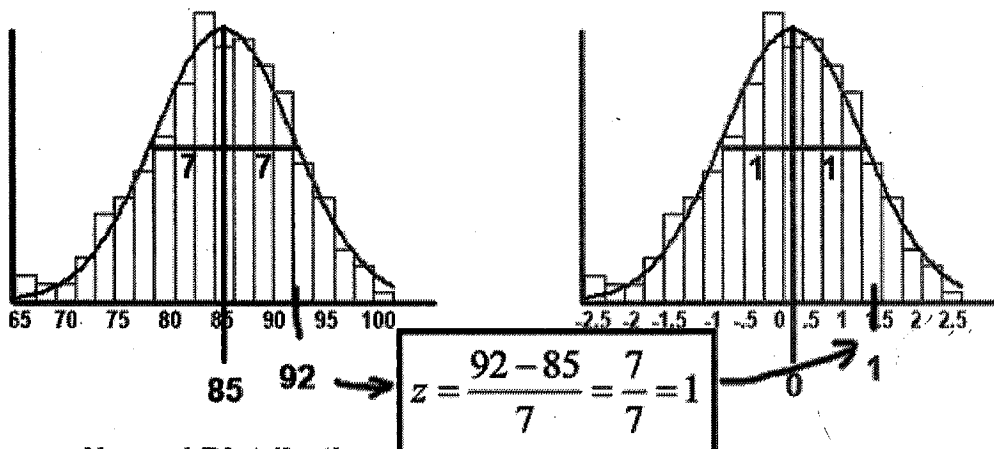
These are called *statistics*.

**Notation for Normal distribution model:  $N(\mu, \sigma)$**

Every time you work a problem involving a Normal distribution, you should sketch the distribution and mark the mean +/- 1 standard deviation, like this:



We can represent the Normal model data values using the original data values (x) or we can standardize by transforming all the x data values into corresponding z-scores. This produces a **Standard Normal Distribution**.



**Normal Distribution**  
(x values)  
 $N(85, 7)$

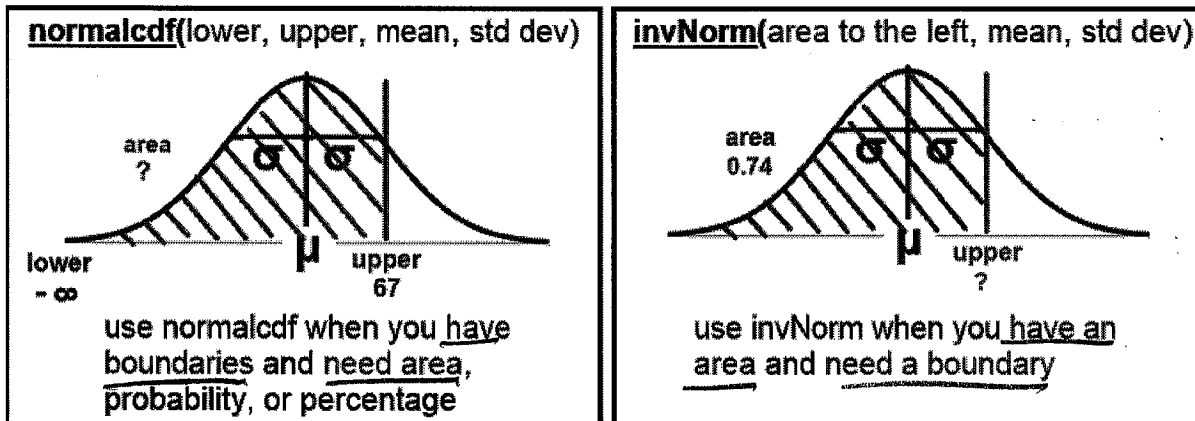
**Standard Normal Distribution**  
(z values)  
 $N(0, 1)$

## Normal distribution model

The Normal distribution model is a **probability density function**, which means:

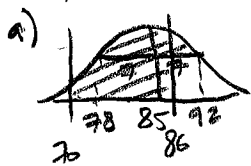
- The area under the whole Normal curve is 1.  
(The sum of all probabilities/percentages = 100%)
- The area under the curve between two x or z values is the percentage of the data in this range (or the probability that one value selected at random will be in this range).

**Two calculator functions for normal distributions:**

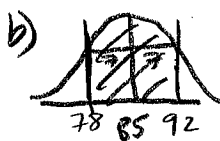


Example: The scores on a test are Normally distributed, with a mean of 85 and a standard deviation of 7.

- What percentage of scores are between 70 and 86?
- What percentage of scores are between 78 and 92?
- What percentage of scores are between  $z=-1$  and  $z=1$ ?
- What percentage of scores are below  $z=-1.5$ ?
- What score is required to be in the top 1%?  
(work both ways, using z then x)



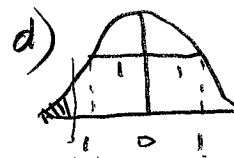
$$P(70 < x < 86) = \text{normalcdf}(70, 86, 85, 7) = .5407 \quad (54.07\%)$$



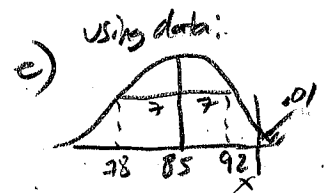
$$P(78 < x < 92) = \text{normalcdf}(78, 92, 85, 7) = .6827 \quad (68.27\%)$$



$$P(-1 < z < 1) = \text{normalcdf}(-1, 1, 0, 1) = .6827 \quad (68.27\%)$$



$$P(z < -1.5) = \text{normalcdf}(-999, -1.5, 0, 1) = .0668 \quad (6.68\%)$$



$$x = \text{invNorm}(.99, 85, 7) = 101.3$$

using z-score:



$$z = \text{invNorm}(.99, 0, 1) = 2.3263$$

$$z = \frac{x - \mu}{\sigma} \Rightarrow 2.3263 = \frac{x - 85}{7}$$

$$\text{solve for } x = 101.3$$

Example: A data set is Normally distributed.

- What percentage of the data is within 1 standard deviation of the mean?
- What percentage of the data is within 2 standard deviations of the mean?
- What percentage of the data is within 3 standard deviations of the mean?

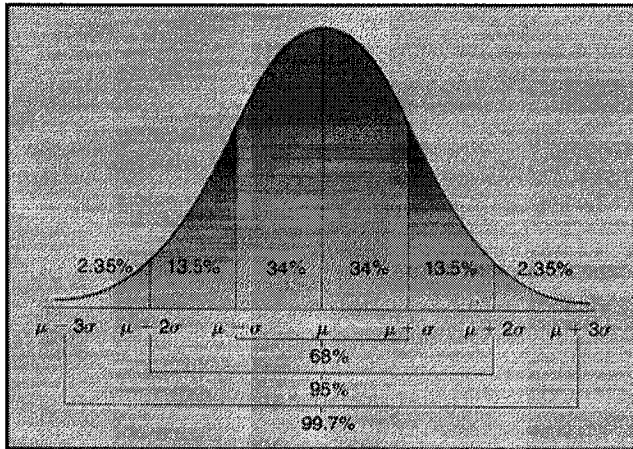
a)  $\text{normalcdf}(-1, 1, 0, 1) = .683$

b)  $\text{normalcdf}(-2, 2, 0, 1) = .954$

c)  $\text{normalcdf}(-3, 3, 0, 1) = .997$

**The 68-95-99.7 Rule...**

*(memorize the 68-95-99.7 % values)*



What percentage of data is above  $2\sigma$  ?

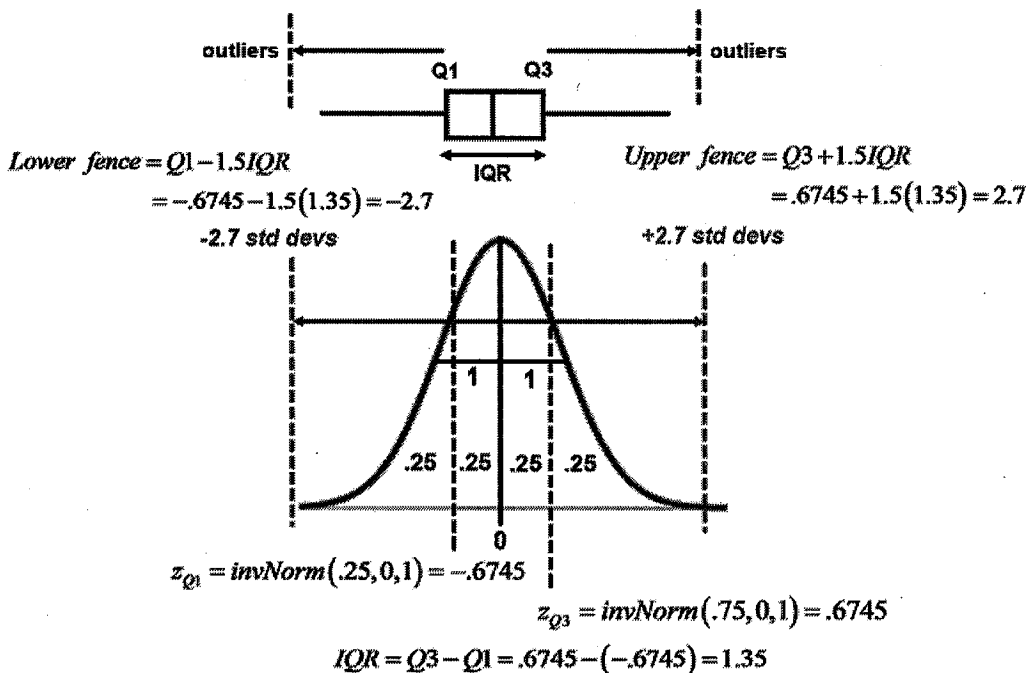
We still haven't answered this question...

"How many standard deviations away from the mean does a data value have to be to be considered 'significantly different' than the mean?"

**$2\sigma$  from the mean**  
= top/bottom 2.5%

**$3\sigma$  from the mean**  
= top/bottom 0.15%

How far away is 'unusual'?



For normal distributions...

- +/- 2 standard deviations from the mean is the top/bottom 2.5% (for normal distributions).
- The IQR rule for outliers ('fences') are at +/- 2.7 standard deviations from the mean (for normal distributions).
- +/- 3 standard deviations from the mean is the top/bottom 0.15% (for normal distributions).

Most people would say 'unusual' starts somewhere between 2 and 3 standard deviations away from the mean.

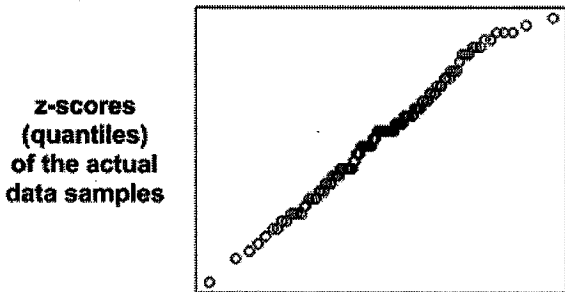
## Determining if it is appropriate to use a Normal distribution

In earlier chapters, we've said that we should only use mean and standard deviation if the distribution is **unimodal and symmetrical**. This is referred to as **Nearly Normal Condition**.

Two ways to determine if a dataset is nearly Normal condition:

- Construct a histogram.
- Construct a Normal Probability Plot (NPP), also known as a Quantile Plot.

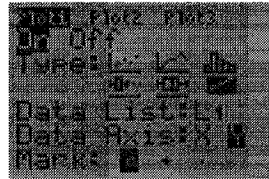
## Normal Probability (Quantile) Plots



z-scores (quantiles) of the samples if they followed a perfect, theoretical Normal distribution

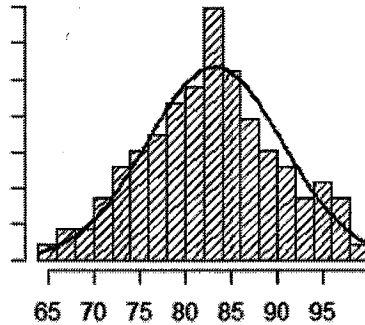
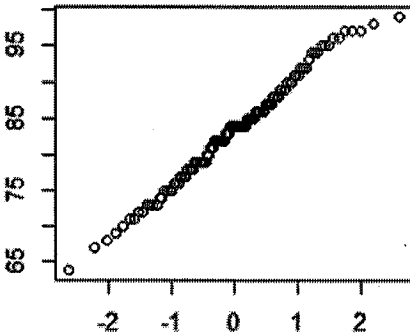
### Use calculator:

- 1) Enter data into list L1.
- 2) Clear Y= equations.
- 3) 2nd Y= (Stat Plot):



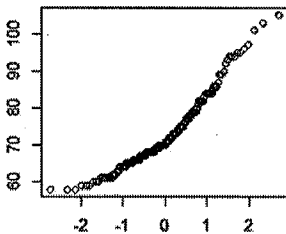
- 4) Zoom 9: ZoomStat

## Interpreting Normal Probability (Quantile) Plots

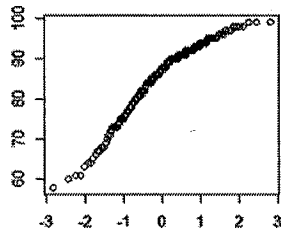


If the NPP is close to a straight line, then the data is close to a Normal distribution shape.

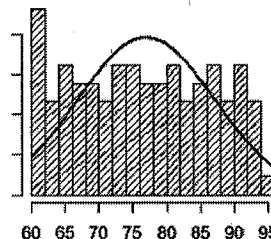
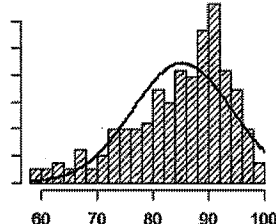
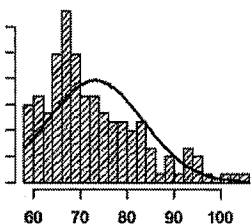
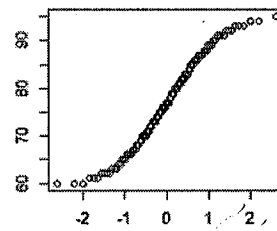
### Skewed Right



### Skewed Left

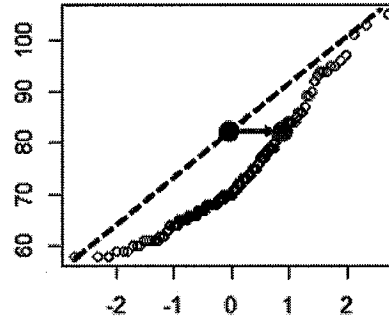
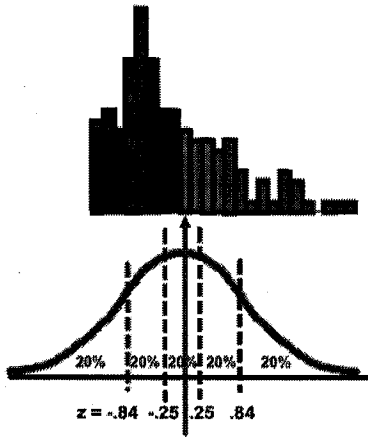


### Uniform



## How does a Normal Probability Plot work?

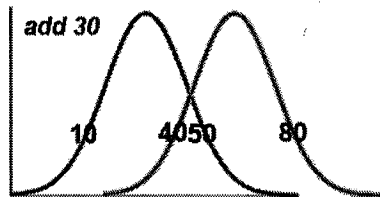
For a perfect, theoretical, Normal distribution, we could use `invNorm` to calculate z-scores for the boundaries between each group of 20% of the population.



Actual data at location marked with red arrow is in the 4th quintile but in 3rd quintile in model so actual z-score will be higher than theoretical. All points in middle of data set are shifted to the right of theoretical, but catch up to line at end of distribution - most of 5th (gray) quintile is in 5th quintile of model as well. This results in the curved plot.

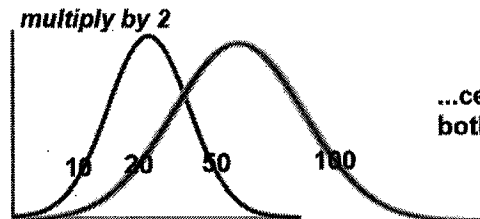
## How does a distribution change if we add or multiply all data by a constant?

Add a constant to each data value...



...center changes, but spread does not change.

Multiply each data value by a constant...



...center and spread both change.

If  $y = ax + b$  (multiply by  $a$ , add  $b$  to each data value)

$$\mu_y = a\mu_x + b$$

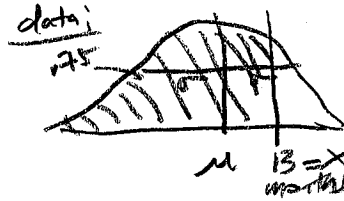
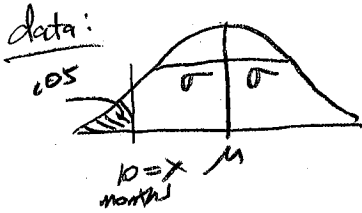
$$med_y = a med_x + b$$

$$\sigma_y^2 = a^2 \sigma_x^2$$

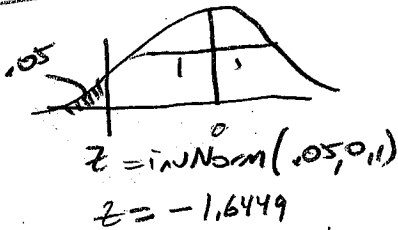
$$\sigma_y = |a| \sigma_x$$

$$IQR_y = |a| IQR_x$$

47. **First steps.** While only 5% of babies have learned to walk by the age of 10 months, 75% are walking by 13 months of age. If the age at which babies develop the ability to walk can be described by a Normal model, find the parameters (mean and standard deviation).



Z-scores



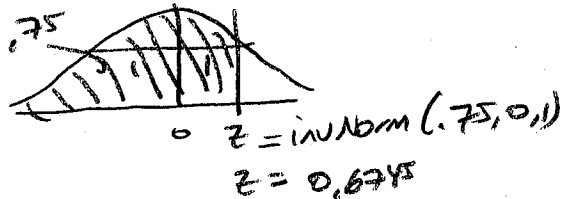
$$z = \frac{x - \mu}{\sigma}$$

$$-1.6449 = \frac{10 - \mu}{\sigma}$$

$$-1.6449\sigma = 10 - \mu$$

$$\mu - 1.6449\sigma = 10$$

Z-scores



$$0.6745 = \frac{13 - \mu}{\sigma}$$

$$0.6745\sigma = 13 - \mu$$

$$\mu + 0.6745\sigma = 13$$

a system of equations:

$$\begin{cases} \mu - 1.6449\sigma = 10 \\ \mu + 0.6745\sigma = 13 \end{cases}$$

fastest way to solve: augmented matrix / RREF:

$$\left[ \begin{array}{cc|c} 1 & -1.6449 & 10 \\ 1 & 0.6745 & 13 \end{array} \right] \rightarrow \text{rref} \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 12.1276 \\ 0 & 1 & 1.29 \end{array} \right]$$

$$\boxed{\begin{aligned} \mu &= 12.1276 \text{ months} \\ \sigma &= 1.29 \text{ months} \end{aligned}}$$