

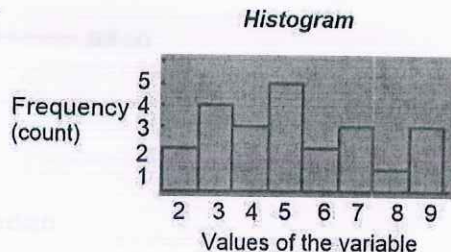
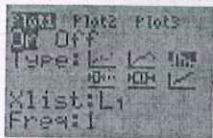
Honors Finite Mathematics – Lesson Notes: Unit 7 (Chapter 9)

9.4 – Statistics: Introduction, Measures of Central Tendency (Center)

Given a set of data, it is always helpful to view the distribution (histogram):

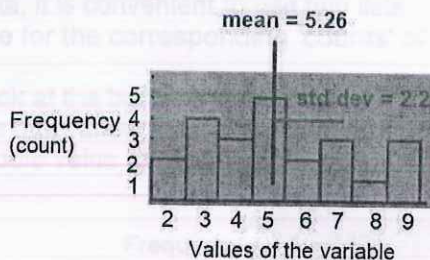
2, 2, 3, 3, 3, 3, 4, 4, 4, 5, 5, 5, 5, 5, 6, 6, 7, 7, 7, 8, 9, 9, 9

- 1) 2nd, '+' (Mem), ClearAllLists, <enter>
- 2) STAT, Edit
- 3) Enter the data set into list 1 (L1)
- 4) 2nd, Y=, set as follows:



- 4) 2nd, Mode to exit
- 5) y=, clear out any equations
- 6) Zoom, 9:StatPlot
- 7) Window, Xscl=1
- 8) Graph

But it is also useful to have a single number which summarizes something about a data distribution. Such a number is called a **statistic**.



Some numbers, such as 'mean' are measures of central tendency (or center)

Some numbers, such as 'standard deviation' are measures of dispersion (or spread)

Measures of central tendency (center):

Mean

$$\bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Median

- Arrange data in order.
- Median is middle value.
- If n is even, median is average of two middle values.

Mode

- Most frequently occurring value or values.
- Can be more than one mode (bimodal, etc.)

If you have a complete data set and can enter it into a calculator, the calculators '1-Var Stats' function can provide mean and median:

- 1) Enter data into L1
- 2) Stats, right arrow to 'CALC'
- 3) 1-Var Stats

```
1-Var Stats
List:L1
FreqList:
Calculate
```

Result...

```
1-Var Stats
x̄=5.260869565
Σx=121
Σx²=743
Sx=2.199532829
σx=2.151185544
n=23
```

← mean

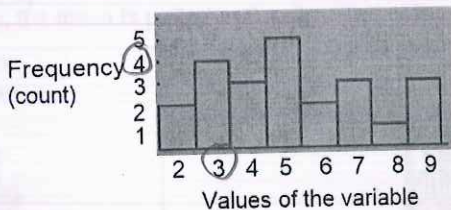
```
minX=2
Q1=3
Med=5
Q3=7
maxX=9
```

← median

Grouped data: Sometimes, we don't have a complete data set, but we have information on groups of data. If we were just given the histogram, we know that there are 4 data values somewhere between 2.5 and 3.5, but we don't know exactly what those values are. They might be 3, 3, 3, 3 or they might be 2.7, 2.8, 3.1, 3.2.

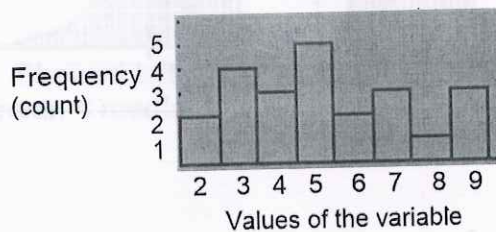
Even with grouped data, we can obtain estimated statistics. To enter the data, it is convenient to use two lists: one for the data values (L1) and one for the corresponding 'counts' of the data values.

Look at the bar for the data value 3. There is a count of 4 data values. We don't know their exact values, but we will assume they are all at the middle value for this bar which is '3'.



So we can enter this histogram's grouped data into lists L1, and L2:

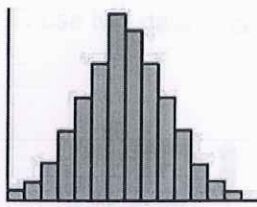
| data values | | counts |
|-------------|----|--------|
| L1 | L2 | |
| 2 | 2 | |
| 3 | 4 | |
| 4 | 3 | |
| 5 | 5 | |
| 6 | 2 | |
| 7 | 3 | |
| 8 | 1 | |
| 9 | 3 | |



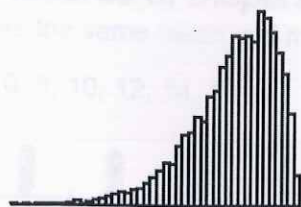
- 1) Enter data into L1 and counts into L2
- 2) Stats, right arrow to 'CALC'
- 3) 1-Var Stats

```
1-Var Stats
List:L1
FreqList:L2
Calculate
```

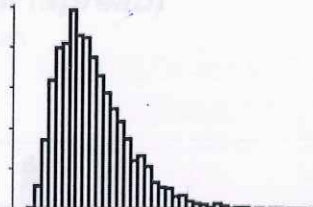
Shape: Symmetric or skewed



Symmetrical



Skewed left
(tail toward lower values)



Skewed right
(tail toward higher values)

How mean, median, and skew are related

What is the mean and median of this data set?

6 7 8 9 10

$$\bar{X} = \frac{6+7+8+9+10}{5} = 8$$

Median = 8

If we move the '10' farther away from the mean, how does mean and median change?

6 7 8 9 20

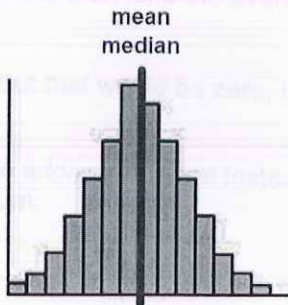
$$\bar{X} = \frac{6+7+8+9+20}{5} = 10$$

Median = 8

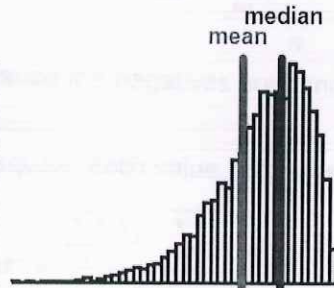
The median is unaffected, but the outlier 'pulls the mean towards itself'

The median is unaffected, but the outlier 'pulls the mean towards itself'

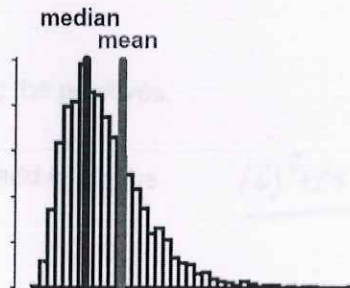
In a perfectly symmetric distribution, the mean and median are aligned, but in skewed distributions, the mean is pulled in the direction of the tail:



Symmetrical



Skewed left
(mean < median)

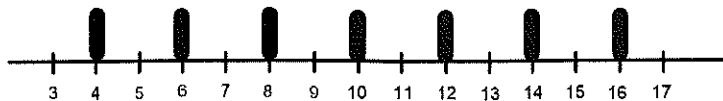


Skewed right
(mean > median)

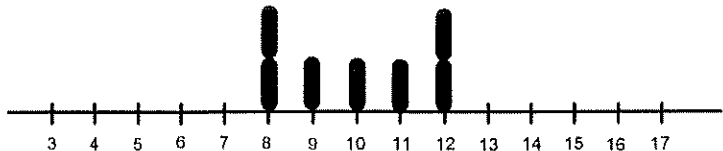
9.5 – Statistics: Measures of Dispersion (spread)

These two data sets have the same mean and median...

4, 6, 8, 10, 12, 14, 16

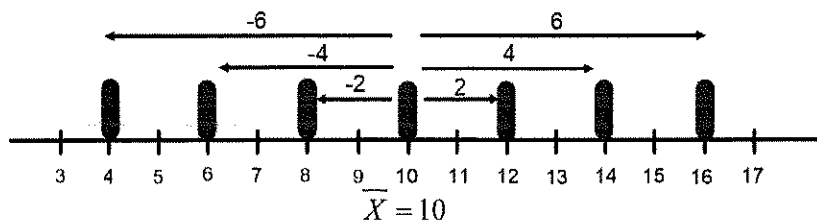


8, 8, 9, 10, 11, 12, 12



...but the data in the top set is spread out further from the mean.

We might ask, 'what is the typical, or average, distance data points are away from the mean?'



$$\frac{6+4+2-6-4-2}{6} = 0$$

For each point, we could compute the difference between the data value and the mean... $X - \bar{X}$

...and then take the average of these values: $\frac{\sum_{i=1}^n (X_i - \bar{X})}{n}$

But that would be zero, because the negatives are canceling the positives.

To address this, we instead square each value before as we add it into the sum:

$$\text{variance} = \sigma^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$$

$$\frac{(6)^2 + (4)^2 + (2)^2 + (-6)^2 + (-4)^2 + (-2)^2}{6}$$

$$\sigma^2 = 18,66$$

The result is the **variance** which is a statistic that gives the average distance squared that the data points are spread around the mean.

There are two problems with variance as a measure of the spread:

- 1) It isn't 'intuitive' to think in terms of 'distances squared'
- 2) The variance is in terms of units squared. For example, if this data represented lengths of fish in cm, the variance would be in cm^2

To get a better measure of dispersion (or spread) that is more intuitive and matches the units of the data set, we just take the square-root of the variance. This is called **standard deviation**:

$$\text{standard deviation} = \sigma = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}}$$

Standard deviation represents the average distance that the data points are from the mean.

Measures of Dispersion

There are many different ways to represent dispersion. A few of the more commonly used are:

1. **Range**-- the difference between the largest value and the smallest value.
2. **Variance**-- deviation from the mean (averaged square deviation)

$$\sigma^2 = \frac{(x_1 - \bar{X})^2 + (x_2 - \bar{X})^2 + \dots + (x_n - \bar{X})^2}{n} = \sum \frac{(x_i - \bar{X})^2}{n}$$

3. **Standard Deviation**-- square root of the variance, 'the average distance that data points are from the mean':

$$\sigma = \sqrt{\frac{(x_1 - \bar{X})^2 + (x_2 - \bar{X})^2 + \dots + (x_n - \bar{X})^2}{n}} = \sqrt{\sum \frac{(x_i - \bar{X})^2}{n}}$$

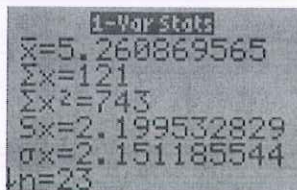
In general, a relatively small standard deviation indicates that the measures tend to cluster close to the mean.

Relatively high standard deviation shows that the measures are widely scattered from the mean.

Using a calculator to compute standard deviation:

- 1) Enter the data either into L1, or for grouped data, into L1 and counts in L2.
- 2) STAT, right arrow, CALC, 1-Var Stats

Results...



```
1-Var stats
x̄=5.260869565
Σx=121
Σx²=743
sx=2.199532829
σx=2.151185544
n=23
```

The calculator gives two different versions of standard deviation:

← s
← σ

The difference between s and σ is subtle. Usually, you should use the s value.

Populations vs. samples from populations

Usually, when we are looking at data that has been collected, it has been collected as being a representative **sample** from a large population about which we are trying to draw conclusions.

For example, let's say we are looking at some data measured using the students in this room. We can calculate statistics about this data, but we would usually be measuring this room's students because we are trying to learn something about **all** finite math students. That would make 'all finite math students' the population, and this particular class' students a **sample** from this population.

If your data is a sample - use s

If your data is the entire population - use σ

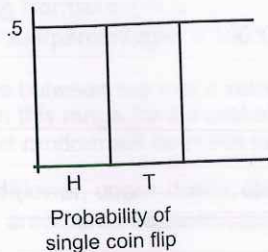
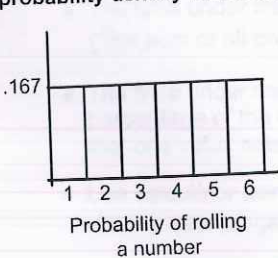
The difference between s and σ is that in the s version we actually divide by $n-1$, not n :

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}} \qquad s = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}}$$

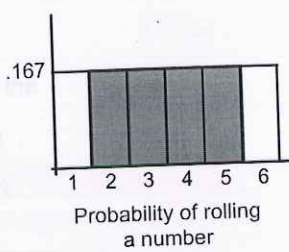
If you want to know more about this difference, please read the posted explanation 'Why is it $n-1$ instead of n ?' on www.mrfelling.com in the 'extras' section of our class page.

9.6: The Normal Distribution and Z-scores

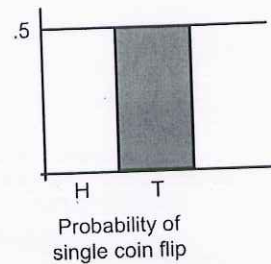
Uniform probability distributions (probability density functions)



Area under probability density curve represents probability of an event:



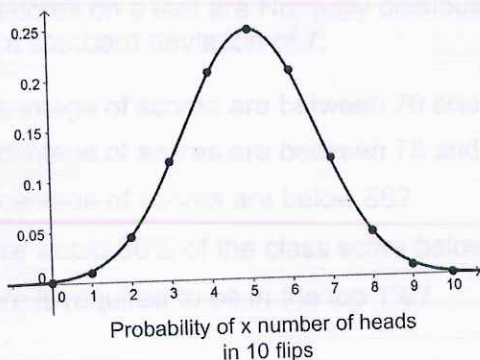
$$P(2,3,4,5) = \frac{4}{6} = \frac{2}{3}$$



$$P(T) = \frac{1}{2}$$

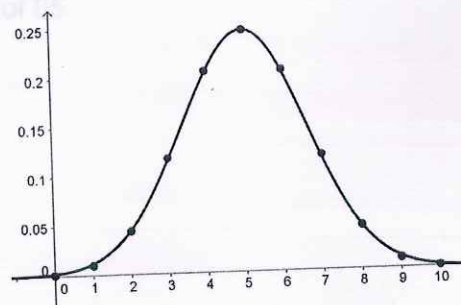
Toss a fair coin 10 times. What is the probability of getting no heads? 1 head? 2 heads?

- $b(10, k; 0.5)$
- $b(10, 0; 0.5) = .0010$
- $b(10, 1; 0.5) = .0098$
- $b(10, 2; 0.5) = .0439$
- $b(10, 3; 0.5) = .1172$
- $b(10, 4; 0.5) = .2051$
- $b(10, 5; 0.5) = .2461$
- $b(10, 6; 0.5) = .2051$
- $b(10, 7; 0.5) = .1172$
- $b(10, 8; 0.5) = .0439$
- $b(10, 9; 0.5) = .0098$
- $b(10, 10; 0.5) = .0010$



Data from many experiments follow this common pattern:

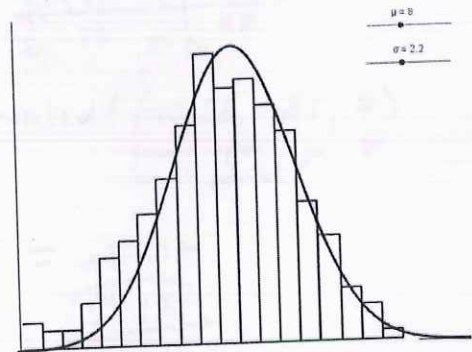
- Heights of adults
- Weights of adults
- Test scores
- Coin tossing



Normal distribution (or Gaussian distribution)

Normal distribution model

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 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')

English alphabet letters for data...

- \bar{x} = mean
- s = standard deviation

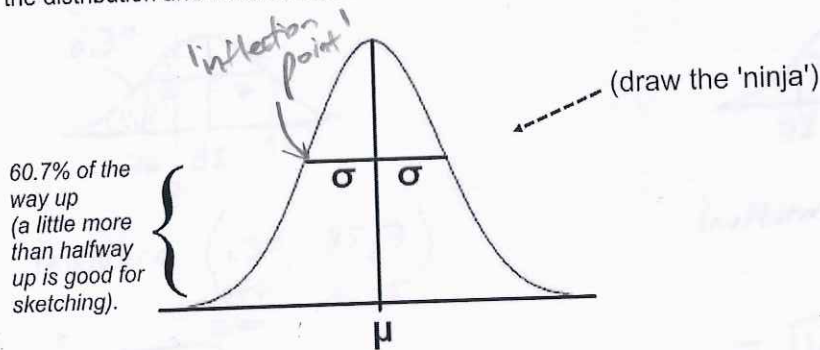
Greek letters for models (populations)...

- μ = mean (mu 'mew')
- σ = standard (sigma) deviation

These are called the **parameters** of the model.

Normal distribution model

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



60.7% of the way up (a little more than halfway up is good for sketching).

Normal distribution model

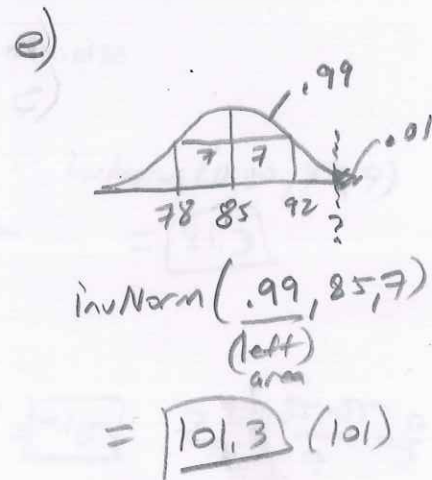
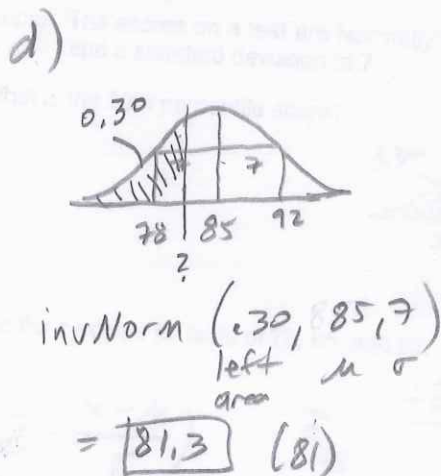
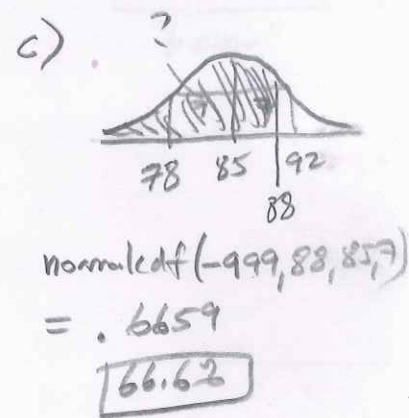
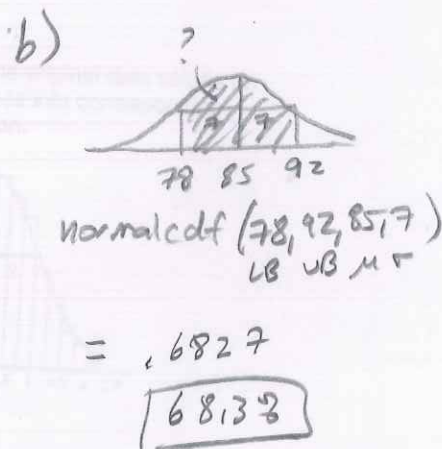
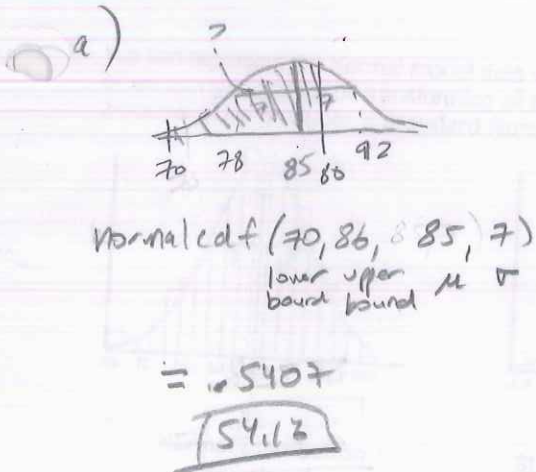
The Normal distribution model is a **probability density function**, which means:

- 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).
- Use calculator **normalcdf**(lower, upper, mean, std dev)
(x or z data range ----> area/percentage/probability)
- Use calculator **invNorm**(area to the **left**, mean, std dev)
(area/percentage/probability ----> x or z data range)

(easiest to see how this works by looking at examples)...

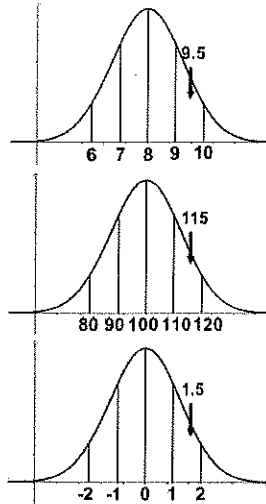
Example: The scores on a test are Normally distributed, with a mean of 85 and a standard deviation of 7.

- What percentage of scores are between 70 and 86?
- What percentage of scores are between 78 and 92?
- What percentage of scores are below 88?
- What score would 30% of the class score below?
- What score is required to be in the top 1%?



Different data set normal curves have different means and standard deviations, but the same shape.

It's useful to standardize by converting values to numbers of standard deviations from the mean.



$$Z = \frac{9.5 - 8}{1} = 1.5$$

$$Z = \frac{115 - 100}{10} = 1.5$$

Z-score

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

x = original data

μ = mean of the data distribution

σ = standard deviation

A Z-score is the number of standard deviations away from the mean this data is.

Positive z-score = above the mean.
Negative z-score = below the mean.

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. We can calculate z-scores to see which is more standard deviations above the mean.

Student 1

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

$$Z = 2.12$$

std dev above the mean

Student 2

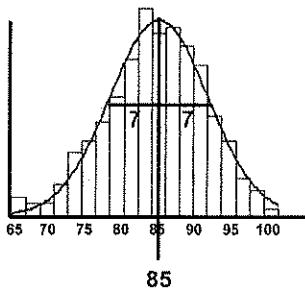
$$Z = \frac{32 - 20.8}{4.8}$$

$$Z = 2.33$$

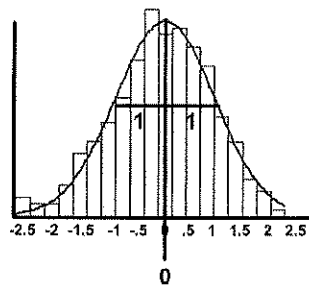
std dev above the mean

"better"

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**.



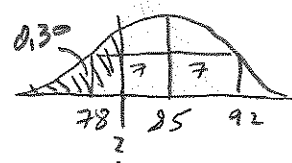
Normal Distribution
(x values)
 $N(85, 7)$



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

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

a) What is the 30th percentile score?



invNorm(0.30, 85, 7)

$$= 81.3$$

b) Find the z-scores for tests of 78, 85, and 92.

$$\left(Z = \frac{x - \mu}{\sigma} \right)$$

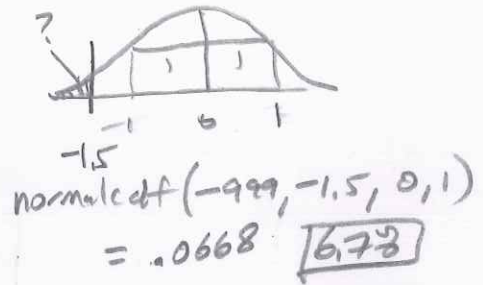
$$z_{78} = \frac{78 - 85}{7} = \frac{-7}{7} = -1\sigma$$

$$z_{85} = \frac{85 - 85}{7} = \frac{0}{7} = 0$$

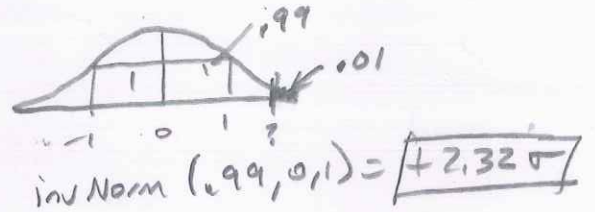
$$z_{92} = \frac{92 - 85}{7} = \frac{7}{7} = +1\sigma$$

c) If a score is 1.5 standard deviations below the mean, what percentage of scores are below this score?

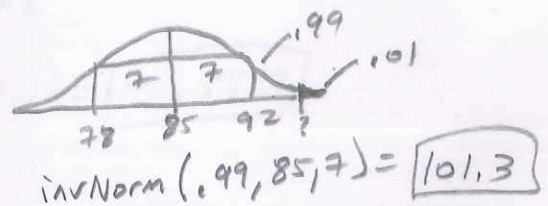
← Z score so use $\mu=0$
 $\sigma=1$
 $z = -1.5$



d) How many standard deviations do you need to be above the mean to be in the top 10%?

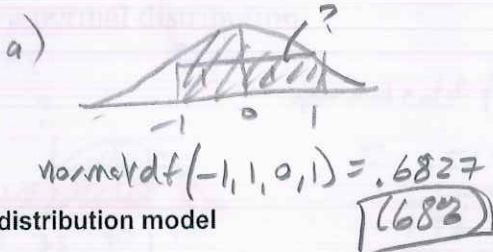


e) What score on the test would put you in the top 10%?

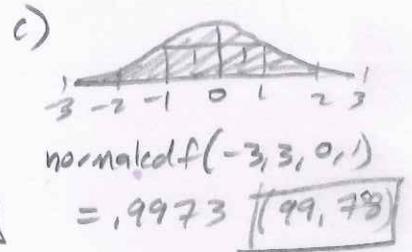
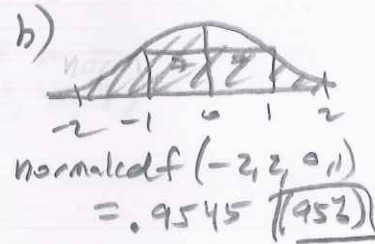


Example: A data set is Normally distributed.

- 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?

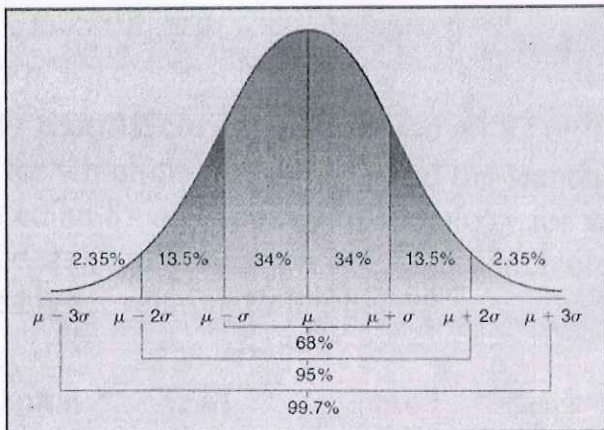


Normal distribution model

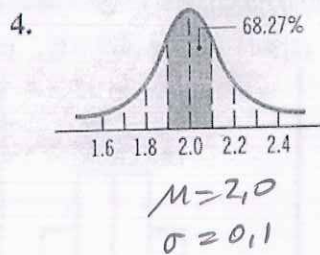
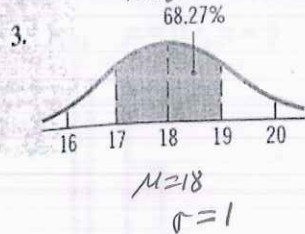
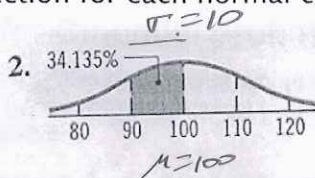
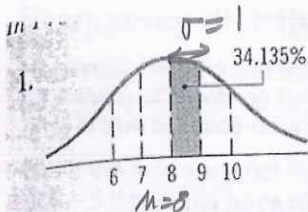


The 68-95-99.7 Rule...

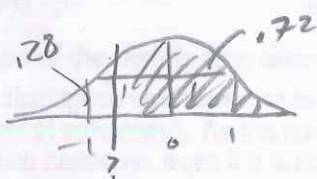
(memorize the 68-95-99.7 % values)



#1-4 Determine μ and σ by inspection for each normal curve.



Find a z-score so that 72% of the population is to the right of z.

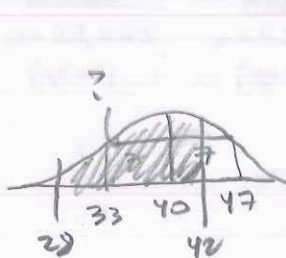


$$\text{invNorm}(0.28, 0, 1) = \boxed{-0.5828}$$

#16 **Life expectancy of clothing** If the average life of a certain make of clothing is 40 months with a standard deviation of 7 months, what percentage of these clothes can be expected to last from 28 months to 42 months? Assume that clothing lifetime follows a normal distribution.

$$\mu = 40 \text{ mo}$$

$$\sigma = 7 \text{ mo}$$



$$\text{normalcdf}(28, 42, 40, 7)$$

$$= 0.5692$$

$$\boxed{56.9\%}$$

Comparing Exam Scores Beth scored an 82 on the final exam in biology for which the mean is 73 and the standard deviation is 9. She scored an 89 on the final in sociology for which the mean is 81 and the standard deviation is 15. In which class was Beth's exam score higher relative to her peers?

$$z_{\text{bio}} = \frac{82-73}{9}$$

$$z_{\text{soc}} = \frac{89-81}{15}$$

$$= 1 \sigma$$

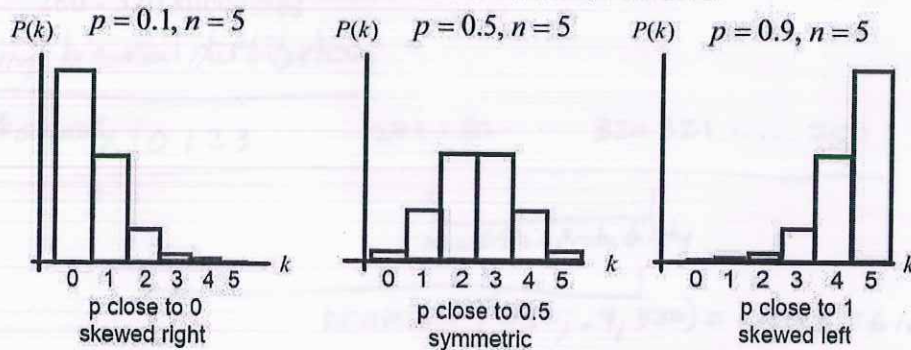
better

$$= 0.67 \sigma$$

Frequency distribution for Binomial probabilities / Bernoulli Trials

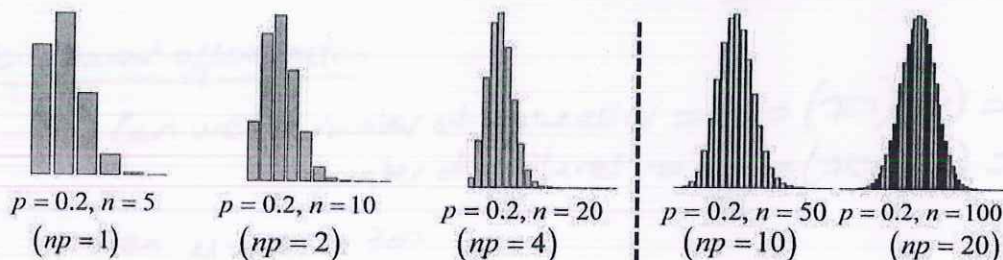
Yesterday, we said that flipping a coin 10 times and determining the probability of receiving 0, 1, 2, etc. heads forms a Normal distribution. That is true because for a fair coin $p=0.5$.

But if the coin were not fair, the Binomial distribution would not be normal: imagine flipping coins 5 times that have different probabilities of landing on heads:



So the shape of the distribution depends upon the probability of success, p .

The Binomial distribution also turns out to depend upon another factor: the sample size (n , the total number of outcomes). As the number of outcomes goes up, a strange and unexpected thing happens: even if p is not close to $.5$, and the distribution is skewed, the distribution becomes more and more symmetrical as n increases, and begins to resemble a Normal shape:



When $np \geq 10$, and $nq \geq 10$ the distribution can be approximated by a Normal distribution.

A Binomial distribution can be approximated by a Normal distribution if:

- p (probability of success) is close to 0.5
- or-
- n is large (at least 10 'successes' and at least 10 'failures', $np > 10$ and $nq > 10$)

For Binomial probabilities that can be approximated by Normal distributions:

$$\mu = np$$

$$\sigma = \sqrt{npq}$$

where n = number of trials

p = probability of success

q = probability of failure ($q = 1 - p$)

24. Suppose a binomial experiment consists of 750 trials and the probability of success for each trial is 0.4. Approximate the probability of obtaining the number of successes indicated by using a normal curve.

280 - 320 successes

Using binomial distribution

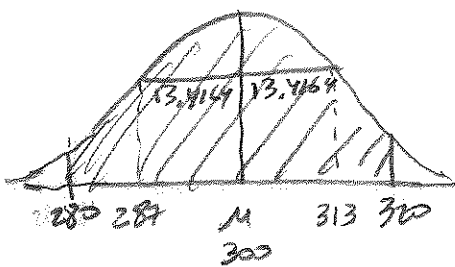
| | | | | | | | | | | | | |
|-------------|---|---|---|---|-----|-----|-----|-----|-----|-----|-----|-----|
| # successes | 0 | 1 | 2 | 3 | ... | 279 | 280 | ... | 320 | 321 | ... | 750 |
| P | $\text{binomcdf}(750, 0.4, 320) = 0.9363616996$ | | | | | | | | | | | |
| | $\text{binomcdf}(750, 0.4, 279) = 0.0628137219$ | | | | | | | | | | | |
| | $= 0.9363616996 - 0.0628137219$ | | | | | | | | | | | |
| | $= \boxed{0.8735479777}$ | | | | | | | | | | | |

Using Normal approximation

Can we? number of successes = $np = (750)(0.4) = 300 \geq 10$ ✓
 number of failures = $nq = (750)(0.6) = 450 \geq 10$ ✓ yes

then $\mu = np = 300$

$\sigma = \sqrt{npq} = \sqrt{750(0.4)(0.6)} = 13.4164$

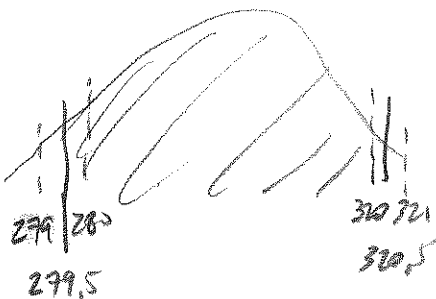


$$P = \text{normalcdf}(280, 320, 300, 13.4164)$$

$$= \boxed{0.86396}$$

(a little low, but close)

Some argue that we should include the part of the distribution between 279-280 and 320-321:



$$P = \text{normalcdf}(279.5, 320.5, 300, 13.4164)$$

$$= \boxed{0.8734826265}$$

$$= \boxed{0.873547}$$
 binomial answer
 (approx accurate to 3 decimal places)

#20 In Math 135 the average final grade is 75.0 and the standard deviation is 10.0. The professor's grade distribution shows that 15 students with grades from 68.0 to 82.0 received Cs. Assuming the grades follow a normal distribution, how many students are in Math 135?

$$\mu = 75.0$$

$$\sigma = 10.0$$



$$\text{normalcdf}(68, 82, 75, 10)$$

$$= .516$$

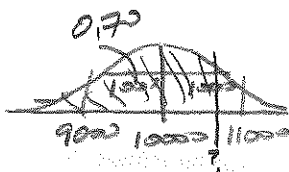
So 51.6% of students should have this grade range. Then:

$$\frac{15}{n} = .51607$$

$$n = \frac{15}{.51607} = 29.066 \quad \boxed{\text{about 29 students}}$$

#18 **Movie Theater Attendance** The attendance over a weekly period of time at a movie theater is normally distributed with a mean of 10,000 and a standard deviation of 1000 persons. Find:

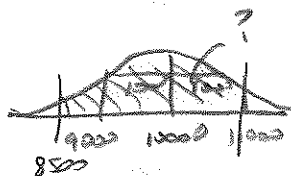
a. the number in the lowest 70% of the attendance figures



$$\text{invNorm}(0.70, 10000, 1000)$$

$$= \boxed{10524 \text{ persons}}$$

b. the percent of attendance figures that falls between 8500 and 11,000 persons.

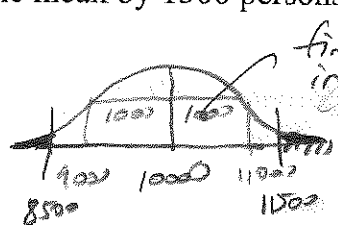


$$\text{normalcdf}(8500, 11000, 10000, 1000)$$

$$= .77454$$

$$\boxed{77.5\%}$$

c. the percent of attendance figures that differs from the mean by 1500 persons or more.



find area in middle = $\text{normalcdf}(8500, 11500, 10000, 1000)$

$$= .8663855...$$

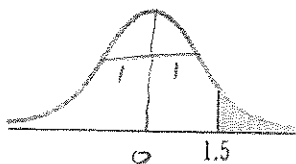
$$\text{shaded} = 1 - .8663855... = .13361$$

$$\boxed{13.4\%}$$

Find the area of the shaded region under the Normal curve:

(no μ, σ given, so assume these are z-scores)

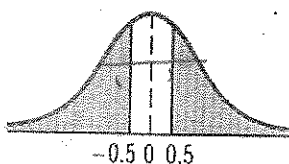
11.



$$= \text{normalcdf}(1.5, 999, 0, 1)$$

$$= \boxed{.0668}$$

12.



$$\text{middle} = \text{normalcdf}(-0.5, 0.5, 0, 1)$$

$$= .38292...$$

$$\text{shaded} = 1 - .38292...$$

$$= \boxed{.617075}$$