

# Stat 107: Introduction to Business and Financial Statistics

## Class 5: Statistics Review, Part II

# The 200 Day SMA System

- Buy when the price closes above the 200 day moving average.
- Sell when the price closes below the 200 day moving average.
- Does this work?

<https://www.portfoliovisualizer.com/test-market-timing-model#analysisResults>

# In the literature a lot



# In the literature

MARKET PULSE

## Apple shares' 200-day moving average ticks lower for first time in nearly two years

Published: Aug 21, 2015 2:46 p.m. ET



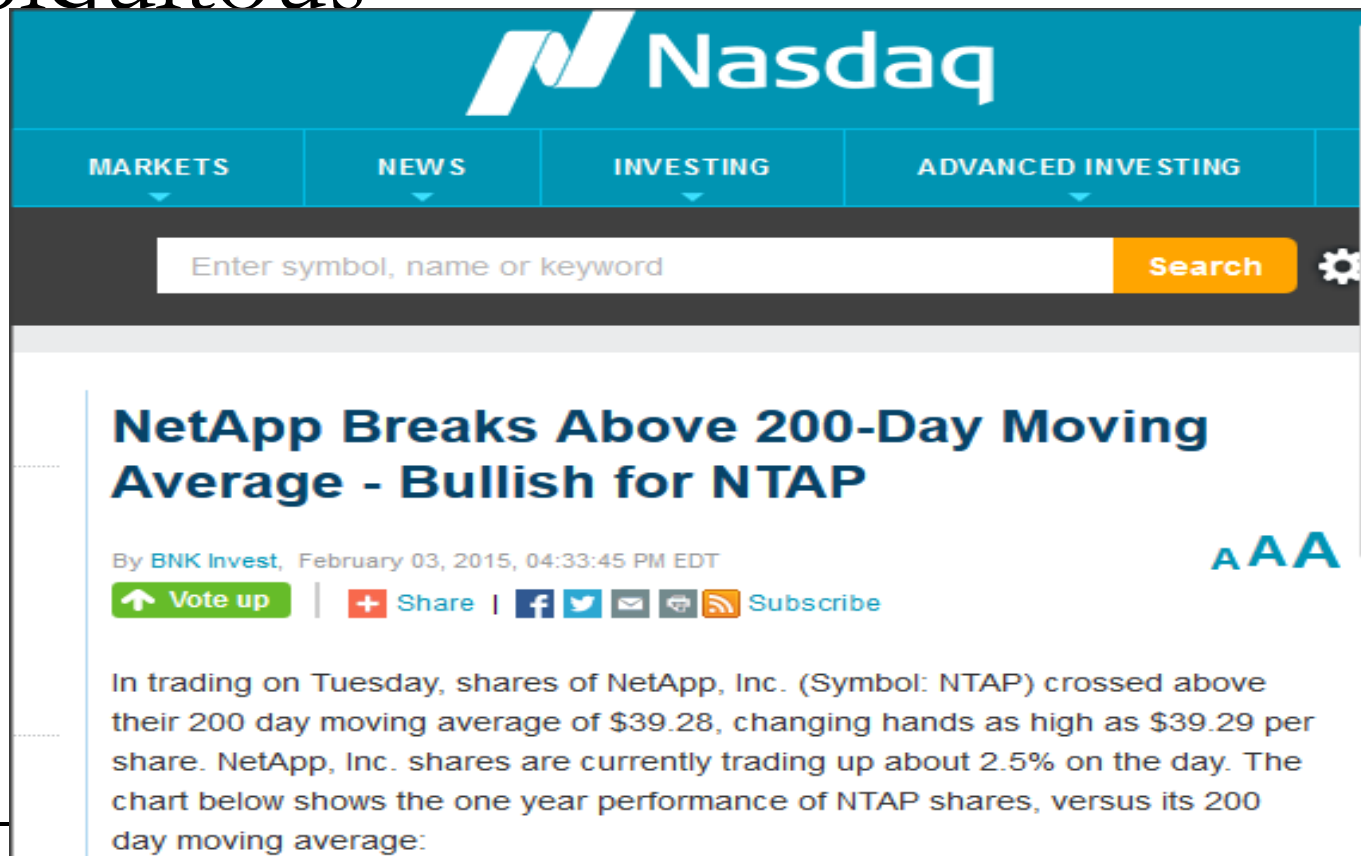
Aa 



By  
**TOMI  
KILGORE**

Apple Inc.'s stock's **AAPL, -1.00%** 200-day moving average is ticking lower in afternoon trade Friday, for the first time in 22 months. This could be significant, because many chart watchers believe bearish technical signals based on the widely-watched long-term moving average are only truly bearish if it is moving lower. For example, many say a "death cross," which is when the 50-day MA crosses below the 200-day MA, isn't a negative signal unless both moving averages are declining. Apple's 200-day MA is currently at \$121.6022, according to FactSet, compared with

# Ubiquitous



The screenshot shows the Nasdaq website interface. At the top is the Nasdaq logo. Below it is a navigation bar with four tabs: MARKETS, NEWS, INVESTING, and ADVANCED INVESTING. A search bar is located below the navigation bar, with the placeholder text "Enter symbol, name or keyword" and a "Search" button. To the right of the search bar is a gear icon. The main content area features a news article titled "NetApp Breaks Above 200-Day Moving Average - Bullish for NTAP". The article is attributed to "BNK Invest" and dated "February 03, 2015, 04:33:45 PM EDT". Below the title is a row of social media sharing buttons: "Vote up", "Share", and "Subscribe". The article text begins with "In trading on Tuesday, shares of NetApp, Inc. (Symbol: NTAP) crossed above their 200 day moving average of \$39.28, changing hands as high as \$39.29 per share. NetApp, Inc. shares are currently trading up about 2.5% on the day. The chart below shows the one year performance of NTAP shares, versus its 200 day moving average:".

**Nasdaq**

MARKETS NEWS INVESTING ADVANCED INVESTING

Enter symbol, name or keyword Search

## NetApp Breaks Above 200-Day Moving Average - Bullish for NTAP

By [BNK Invest](#), February 03, 2015, 04:33:45 PM EDT

[Vote up](#) | [Share](#) | [f](#) [t](#) [e](#) [r](#) [s](#) [Subscribe](#)

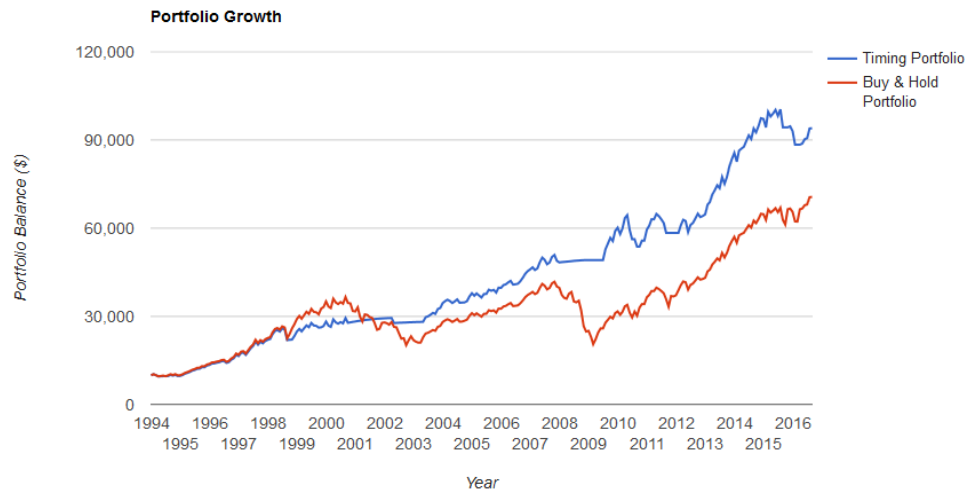
AAA

In trading on Tuesday, shares of NetApp, Inc. (Symbol: NTAP) crossed above their 200 day moving average of \$39.28, changing hands as high as \$39.29 per share. NetApp, Inc. shares are currently trading up about 2.5% on the day. The chart below shows the one year performance of NTAP shares, versus its 200 day moving average:

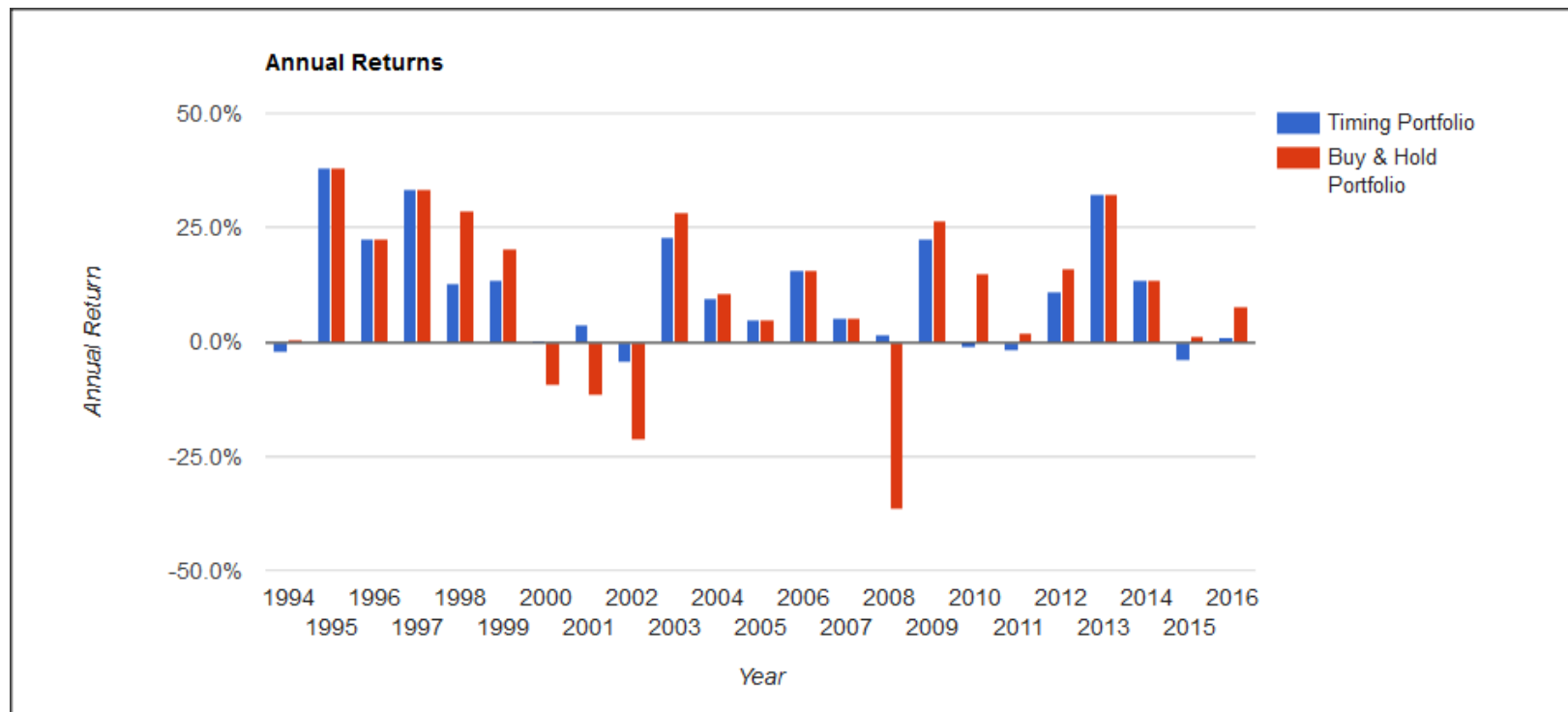
# 200 day MA on SPY

■ What do those flat lines indicate?

Portfolio	Initial Balance	Final Balance	CAGR	Std.Dev.
Timing Portfolio	\$10,000	\$94,024	10.39%	10.36%
Buy & Hold Portfolio	\$10,000	\$70,638	9.01%	14.85%

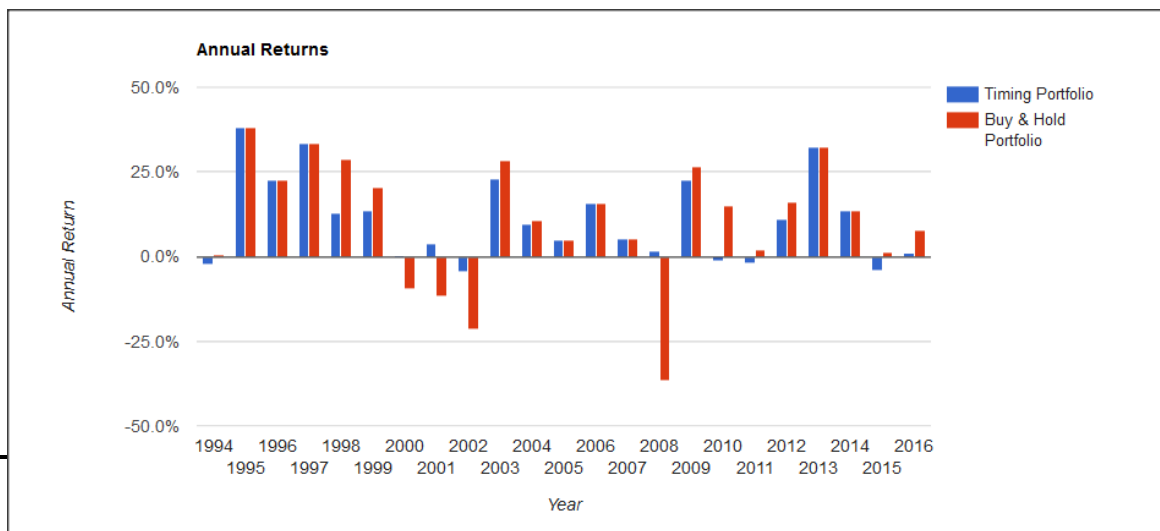


# Look at 2008



# 200 Day MA on AAPL

Portfolio	Initial Balance	Final Balance	CAGR	Std.Dev.
Timing Portfolio	\$10,000	\$304,195	12.20%	37.56%
Buy & Hold Portfolio	\$10,000	\$1,755,624	19.03%	46.61%





# Confusing Return Math

- The annualized return is the geometric average of returns .
- Consider the following annual returns

2003	10.5%
2002	−5.6%
2001	23.4%
2000	−15.7%
1999	8.9%

# Return Statistics

- The cumulative return is

$$(1.105 \times 0.944 \times 1.234 \times 0.843 \times 1.089) - 1 = 18.2\%$$

- The average return is

$$\frac{10.5\% - 5.6\% + 23.4\% - 15.7\% + 8.9\%}{5} = 4.3\%$$

- If this average return is compounded for five years the cumulative return is greater than 18.2%

$$(1.043 \times 1.043 \times 1.043 \times 1.043 \times 1.043) - 1 = 23.4\%$$

# Return Statistics

- Now consider the geometric average return

$$(1.105 \times 0.944 \times 1.234 \times 0.843 \times 1.089)^{\frac{1}{5}} - 1 = 3.4\%$$

- As expected, if the annualized return is compounded for five years, the cumulative return remains the same

$$(1.034 \times 1.034 \times 1.034 \times 1.034 \times 1.034) - 1 = 18.2\%$$

- We will come back to this annualized return concept later in the course.

# Mean > CAGR

Year and Return (%)		Date Range	
2014	13.80	Jan 1	1871 to Dec 31 2014
2013	32.43	<input type="checkbox"/>	Adjust for <a href="#">Inflation</a>
2012	15.88	<input checked="" type="checkbox"/>	Include Dividends
2011	2.07	<input type="button" value="Calculate"/>	
2010	14.87	"Average" return:	10.77 %
2009	27.11	Annualized return	9.11 %
2008	-37.22	(= True CAGR):	18.72 %
2007	5.46	Standard Deviation:	\$ 281,748.60
2006	15.74	\$1.00 grew to:	
2005	4.79		
2004	10.82		
2003	28.72		

# The Normal Distribution

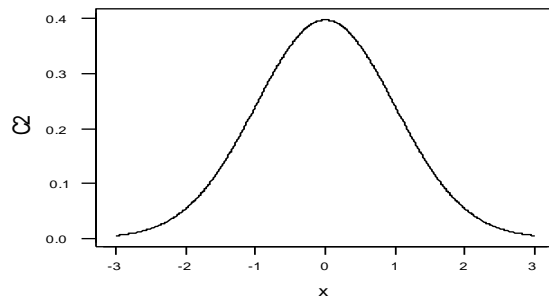
The normal distribution is the most fundamental continuous distribution used in statistics so I guess we better learn it. Many methods that are widely used in economics, finance, and marketing are based on an assumption of a normal distribution.

Normal Dist.

Gaussian Dist.

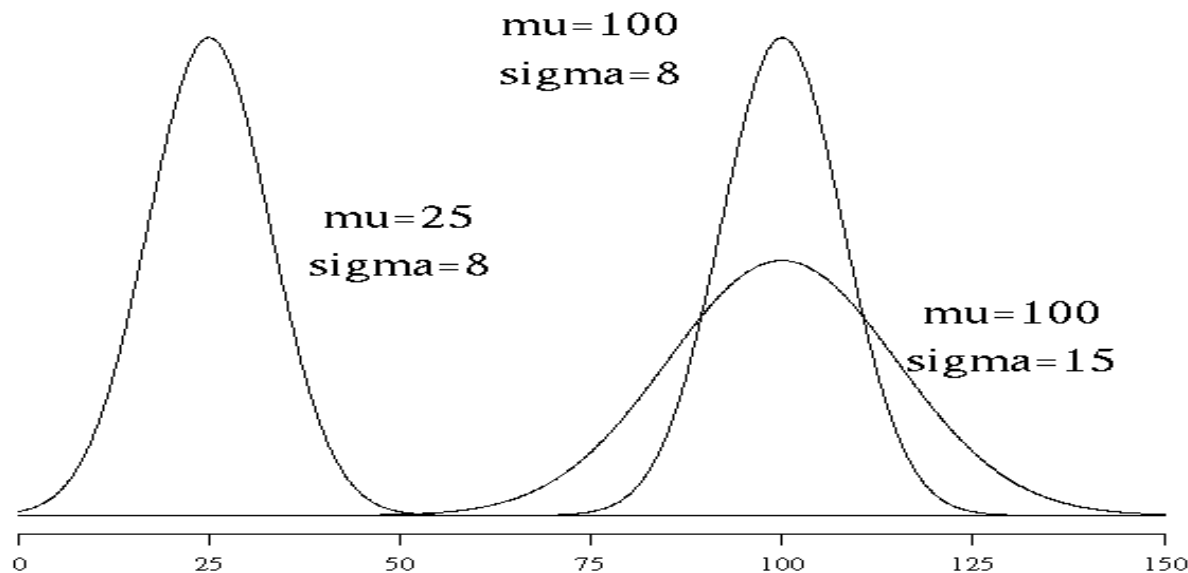
Bell-Shaped Curve

*equivalent names*



# What controls the shape of the curve ?

The normal distribution is governed by the **two parameters** :  $\mu$  (the *mean*) and  $\sigma$  (the *standard devia*

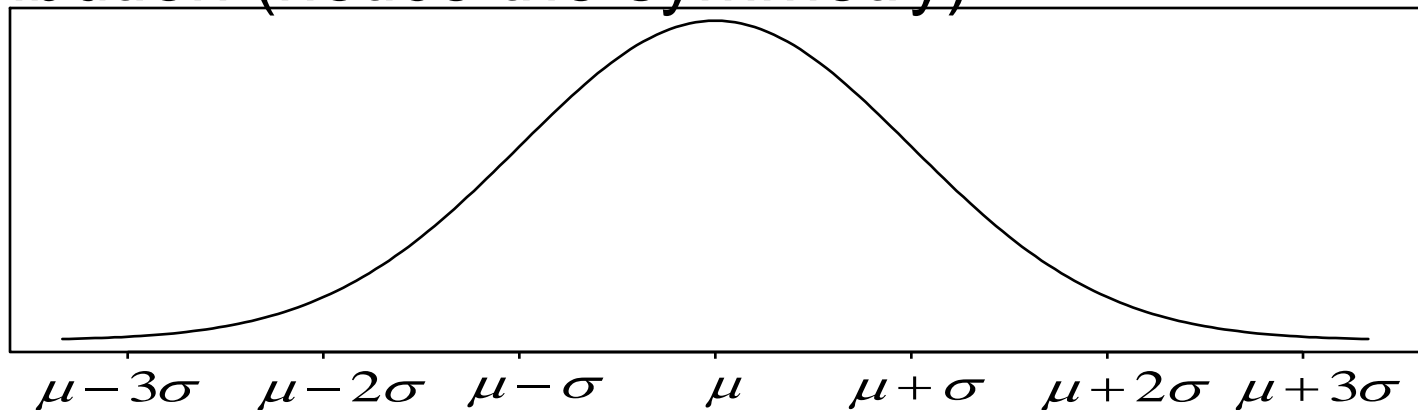


# Notation for the normal curve

We use the notation

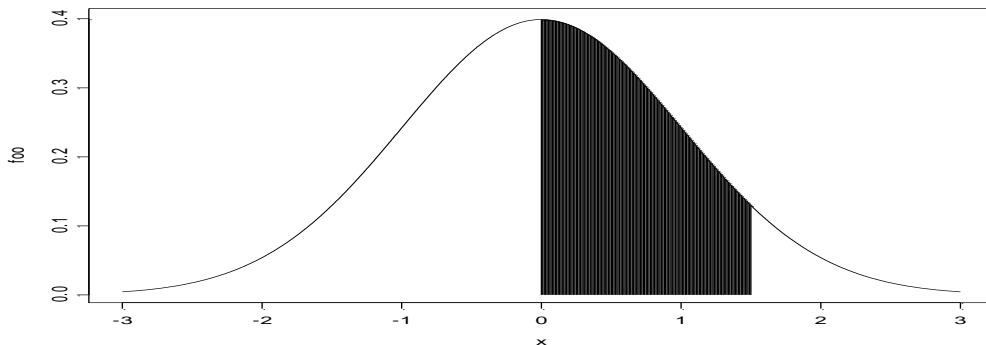
$$X \sim N(\mu, \sigma^2).$$

Here is a general picture of a normal distribution (notice the symmetry).



# Finding normal probabilities

- $P(a \leq X \leq b) = \text{area under the curve between } a \text{ and } b.$



We use the computer or tables to do this



# Using the Computer

## ■ R

□  $P(X < x) = \text{pnorm}(x, m, s)$

□  $f(x) = \text{dnorm}(x, m, s)$

# Example

- Daily stock market returns on the NY Stock Exchange are approximately normally distributed, with a mean of 0.0317% and a standard deviation of 1.046%.
- What percentage of the time are daily returns positive?
- What percentage of the time do daily returns return between 1% and 2%?

# Example

- We are told  $X$  is normal  $\mu = 0.0317$  and  $\sigma = 1.046$
- To find  $P(X > 0)$

```
> 1-pnorm(0, .0317, 1.046)
[1] 0.5120885
```

# Example

- What percentage of the time do daily returns return between 1% and 2%?

```
pnorm(2, .0317, 1.046) - pnorm(1, .0317, 1.046)  
[1] 0.1473609
```

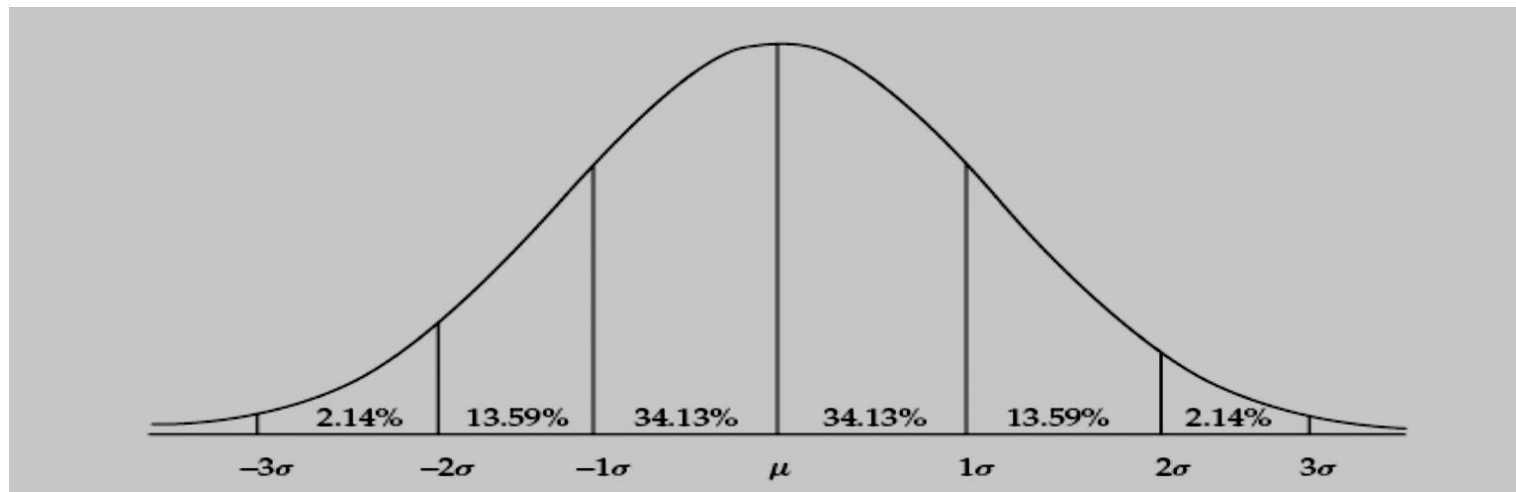
# Normal Coverage Rule

*For  $X \sim N(\mu, \sigma^2)$ ,*

$$P(\mu - \sigma \leq X \leq \mu + \sigma) = 0.68$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 0.95$$

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 0.997$$



# Returns may not be Normal

- It is nice to assume returns are normally distributed but there is an issue
- Simple returns are defined to be

$$R_t = \frac{P_t - P_{t-1}}{P_{t-1}}$$

- As long as the price  $> 0$ ,  $R_t > -1$  [the most you can lose is all your money]

# Returns may not be Normal

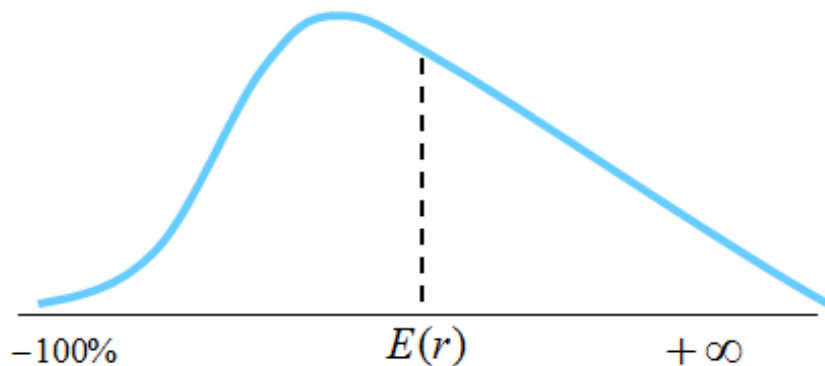
- Suppose returns are normal with mean 0.05 and std dev .5
- Lets calculate  $P(R_t < -1)$

```
pnorm(-1, .05, .5)  
[1] 0.01786442
```

- This implies that there is a 1.8% chance that the asset price will be negative. This is why the normal distribution may not be appropriate for simple returns.

# The Log Normal Distribution

- The Log Normal distribution is sometimes used to compensate for this-more details later on.





# Skewness-measure of symmetry

- Skewness of a random variable is

$$\text{Skew}(X) = E \left[ \left( \frac{X - \mu_X}{\sigma_X} \right)^3 \right]$$

- If  $X$  has a symmetric distribution about  $\mu$  then  $\text{Skew}=0$
- $\text{Skew} > 0 \Rightarrow$  pdf has long right tail, and median  $<$  mean
- $\text{Skew} < 0 \Rightarrow$  pdf has long left tail, and median  $>$  mean

# Skewness Example

- Using the discrete distribution for the return on Microsoft stock shown earlier, we have

$$\begin{aligned}\text{skew}(X) &= [(-0.3 - 0.1)^3 \cdot (0.05) + (0.0 - 0.1)^3 \cdot (0.20) \\ &\quad + (0.1 - 0.1)^3 \cdot (0.5) + (0.2 - 0.1)^3 \cdot (0.2) \\ &\quad + (0.5 - 0.1)^3 \cdot (0.05)] / (0.141)^3 \\ &= 0.0\end{aligned}$$

# Kurtosis-measure heavy tails

- Kurtosis is a measure of whether the data are peaked or flat relative to a normal distribution.
- That is, data sets with high kurtosis tend to have a distinct peak near the mean, decline rather rapidly, and have heavy tails.
- Data sets with low kurtosis tend to have a flat top near the mean rather than a sharp peak.
- Normal data has kurtosis=3 (or 0 depending on routine).

# Kurtosis-measure of tail thickness

- Defined to be

$$\text{Kurt}(X) = E \left[ \left( \frac{X - \mu_X}{\sigma_X} \right)^4 \right]$$

- For our MSFT example

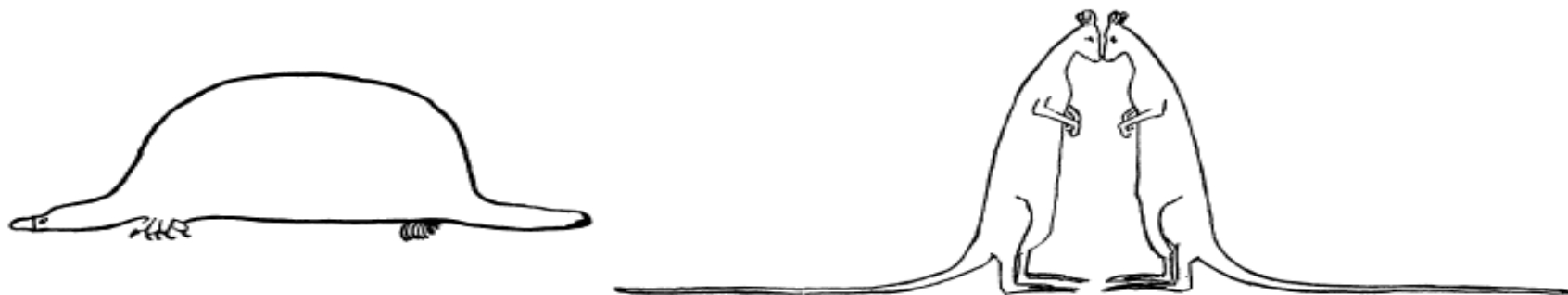
$$\begin{aligned} \text{Kurt}(X) &= [(-0.3 - 0.1)^4 \cdot (0.05) + (0.0 - 0.1)^4 \cdot (0.20) \\ &\quad + (0.1 - 0.1)^4 \cdot (0.5) + (0.2 - 0.1)^4 \cdot (0.2) \\ &\quad + (0.5 - 0.1)^4 \cdot (0.05)] / (0.141)^4 \\ &= 6.5 \end{aligned}$$

# Kurtosis compared to 3

- A normal distribution has a kurtosis of 3
- Kurtosis  $>3 \Rightarrow X$  has fatter tails than normal distribution
- kurtosis  $<3 \Rightarrow X$  has thinner tails than normal distribution

# No Comment

\* In case any of my readers may be unfamiliar with the term “kurtosis” we may define mesokurtic as “having  $\beta_2$  equal to 3,” while platykurtic curves have  $\beta_2 < 3$  and leptokurtic  $> 3$ . The important property which follows from this is that platykurtic curves have shorter “tails” than the



normal curve of error and leptokurtic longer “tails.” I myself bear in mind the meaning of the words by the above *memoria technica*, where the first figure represents platypus, and the second kangaroos, noted for “lepping,” though, perhaps, with equal reason they should be hares!

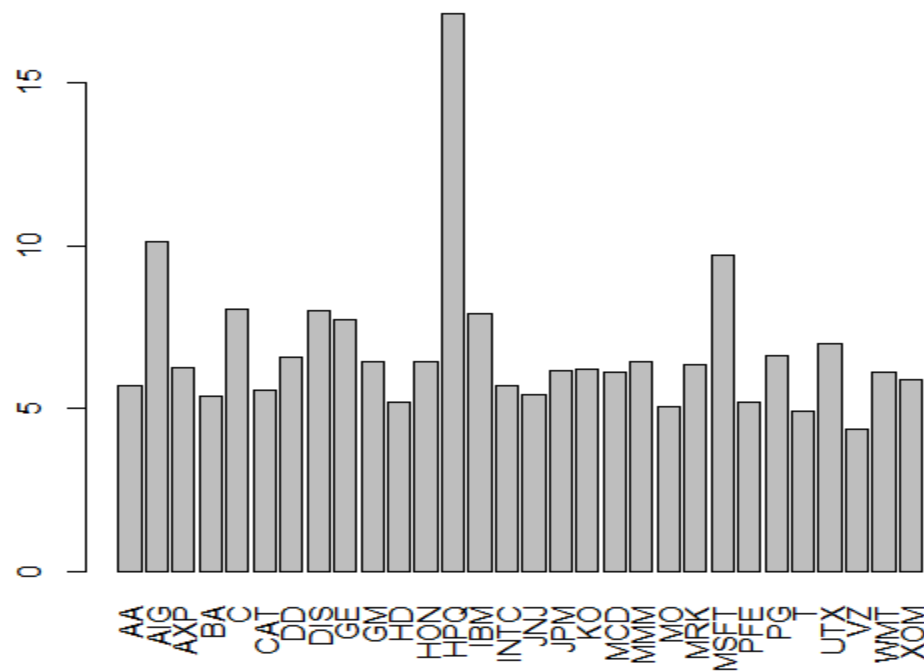
# Calculating Moments in R

■ `install.packages("moments")`

■ `library(moments)`

```
kurtosis(rnorm(1000))  
[1] 2.967884  
> kurtosis(monthlyReturn(Ad(AAPL)))  
monthly.returns  
5.058169  
> kurtosis(monthlyReturn(Ad(PG)))  
monthly.returns  
3.218961  
> kurtosis(dailyReturn(Ad(AAPL)))  
daily.returns  
8.933392  
> kurtosis(dailyReturn(Ad(PG)))  
daily.returns  
10.76613
```

# Kurtosis for the DJ30

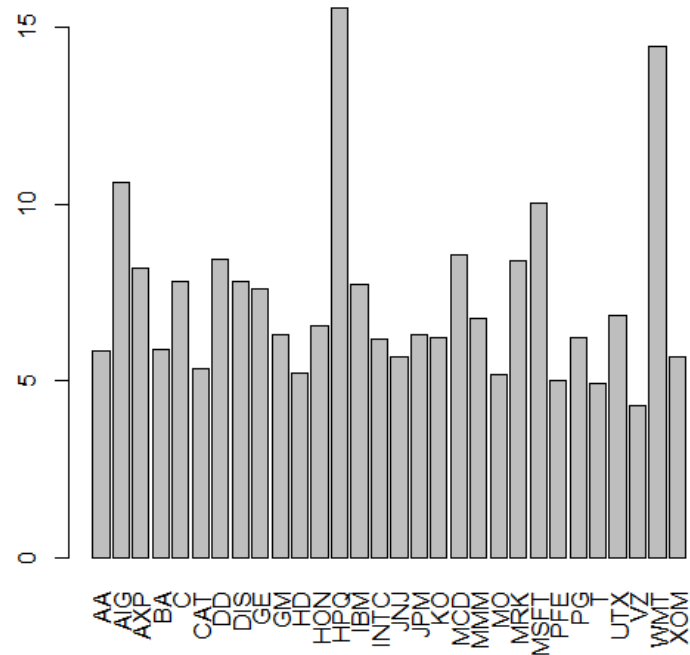




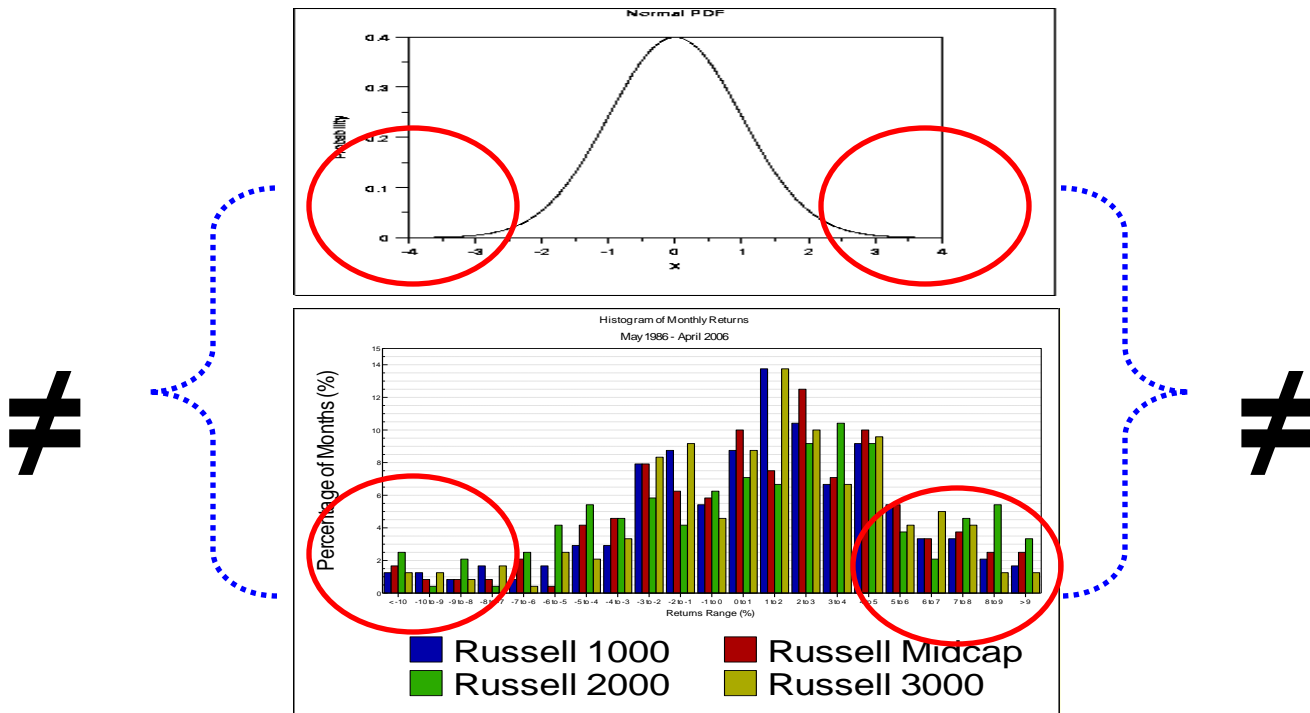
# Code for Previous Graph

```
assets=c("AA","AIG","AXP","BA","C","CAT","DD","DIS","GE","GM","H  
D","HON","HPQ","IBM","INTC","JNJ","JPM","KO","MCD","MMM","MO","M  
RK","MSFT","PFE","PG","T","UTX","VZ","WMT","XOM")  
n=length(assets)  
store.kurts=1:n  
for(i in 1:n) {  
    temp=getSymbols(assets[i],from="2010-01-01",auto.assign=FALSE)  
    store.kurts[i] = kurtosis(dailyReturn(Ad(temp)))  
}  
barplot(store.kurts,las=3)
```

# Kurtosis for the DJ30



# “Deal or No Deal?” Normal or Not Normal?



**Twenty Year Histogram of Monthly Index Returns**

# “If it doesn’t fit, you must acquit.”

- Johnnie Cochran



## The “Fatter” the distribution tails, the less reliable the statistics!

- “If the population of price changes is strictly normal, on average for any stock.....an observation more than five standard deviations from the mean should be observed about once every 7,000 years. In fact such observations seem to occur about once every three or four years” – Eugene Fama, *Journal of Business*, January 1965
- Under the assumption of normal return distributions, the probability of the October 1987 crash was so remote that according to efficient market theory it would have been virtually impossible – Jackwerth and Rubinstein, *Journal of Finance*, Vol 51 1996
- “The problem for traders is that it is much more complicated to create models for a world of fat tails than for a world of bell curves. As a result, traders repeatedly get caught out by “unprecedented” market movements. The collapse of two hedge funds, Long-Term Capital Management in 1998 and Amaranth Advisors in 2006, were cases in point” – *The Economist*, October 18<sup>th</sup> 2007.

# Introduction to Value at Risk

- Consider a  $W_0 = \$10000$  investment in MSFT for one month.
- Assume the return is normally distributed with mean 0.05 and stdev 0.1

$$R \sim N(0.05, (0.1)^2)$$

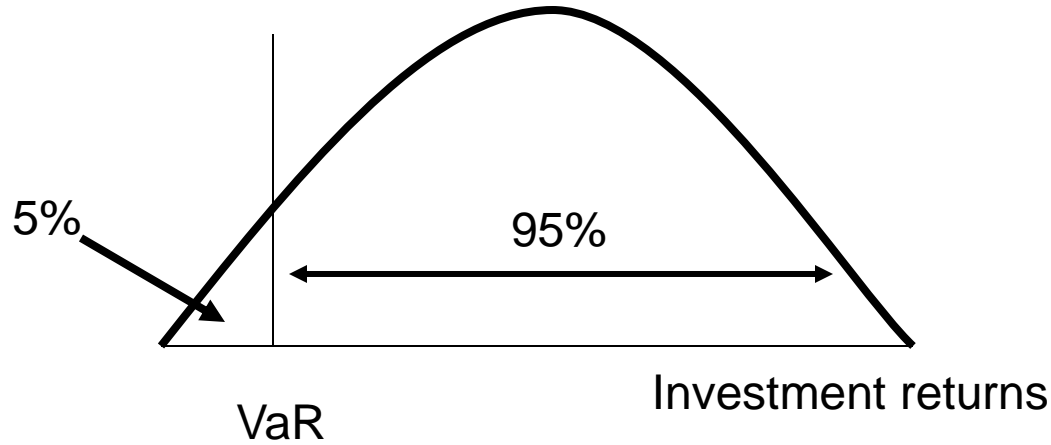
- Goal-calculate how much we can lose with specified probability  $\alpha$

# In the Jargon of VaR

- In the jargon of VaR, suppose that a portfolio manager has a daily VaR equal to \$1 million at 1 percent. This statement means that there is only one chance in 100 that a daily loss bigger than \$1 million occurs under normal market conditions.

Financial Modeling  
Simon Benninga  
third edition

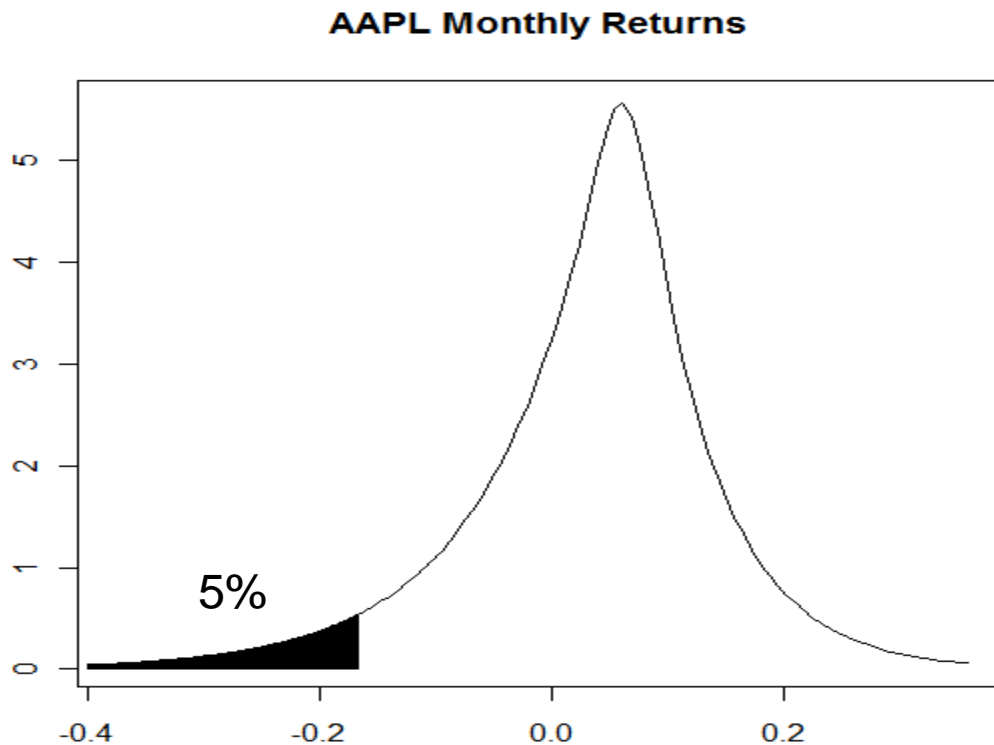
# Choice of confidence level – 95%



Normal market conditions – the returns that account for 95% of the distribution of possible outcomes.

Abnormal market conditions – the returns that account for the other 5% of the possible outcomes.

# Consider AAPL





# Value at Risk

- What is the monthly value-at-risk (VaR) on the \$10,000 investment with 5% probability?

$$R \sim N(0.05, (0.1)^2)$$

```
> qnorm(.05, .05, .1)
[1] -0.1144854
```

- There is a 5% chance of a loss of 11.44% or worse.
- The loss in investment value is  $\$10000(-11.44) = -\$1144$
- Because VAR represents a loss it is usually reported as a positive number.

# Sums of Independent Normals

- If  $X_1$  and  $X_2$  are independent and each normally distributed

$$X_i \sim N(\mu_i, \sigma_i^2)$$

- Then the sum is normally distributed

$$aX_1 + bX_2 \sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$$

- This rule holds for combining any number of independent normal random variables.
-

# Example: Monthly Returns.

- Suppose  $\text{FBALX} \sim N(0.005, .033^2)$
- Suppose  $\text{GSMIX} \sim N(0.004, .015^2)$
- What is  $P(\text{FBALX} > \text{GSMIX})$ ?

# Example

- Let  $D = \text{FBALX-GSMIX}$
- We want to find  $P(D > 0)$ .
- From rule for combining normals we know  
$$D \sim N(.005-.004, .033^2 + .015^2) = N(.001, .036^2)$$

```
> 1-pnorm(0, .001, .0036)
[1] 0.6094085
```

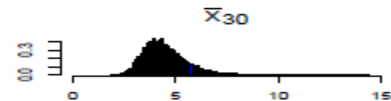
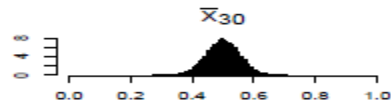
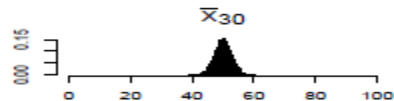
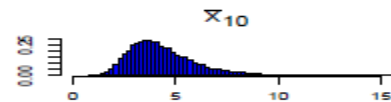
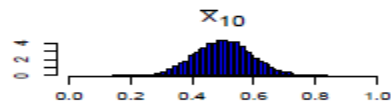
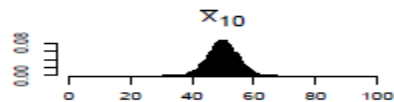
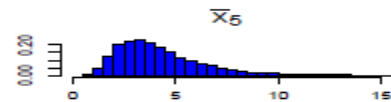
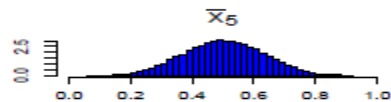
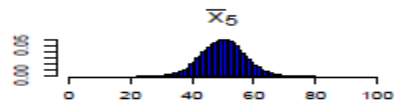
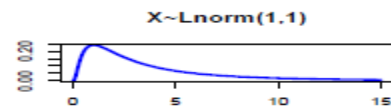
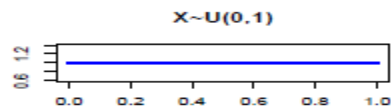
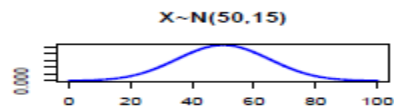
# The Central Limit Theorem

The CLT states that if random samples of size  $n$  are repeatedly drawn from any population with mean  $\mu$  and variance  $\sigma^2$ , then **when  $n$  is large**, the distribution of the sample means will be approximately normal :

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

If the population is normal this is true for any sample size.

# Visually



# Some R Demos

## ■ Nice R library:

---

`clt.examp`

*Plot Examples of the Central Limit Theorem*

---

### **Description**

Takes samples of size `n` from 4 different distributions and plots histograms of the means along with a normal curve with matching mean and standard deviation. Creating the plots for different values of `n` demonstrates the Central Limit Theorem.

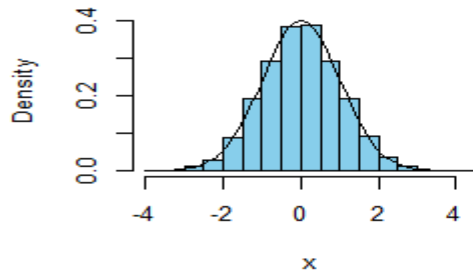
### **Usage**

```
clt.examp(n = 1, reps = 10000, nclass = 16, norm.param=list(mean=0,sd=1),
          gamma.param=list(shape=1, rate=1/3), unif.param=list(min=0,max=1),
          beta.param=list(shape1=0.35, shape2=0.25))
```

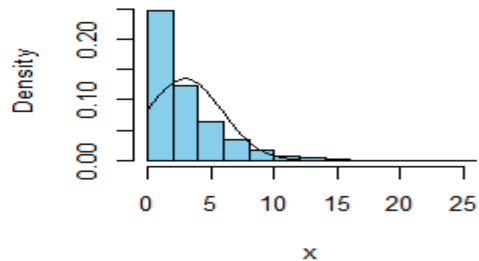
# clt.examp(1)

sample size = 1

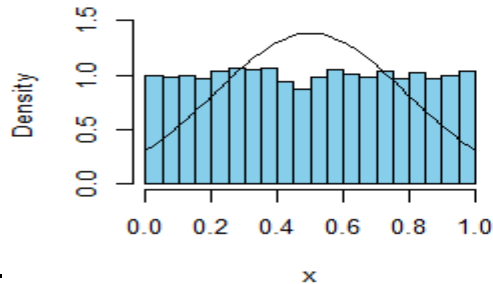
**Normal**



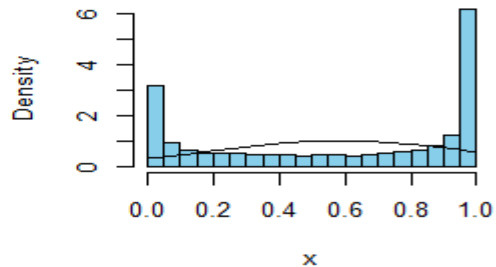
**Gamma**



**Uniform**



**Beta**

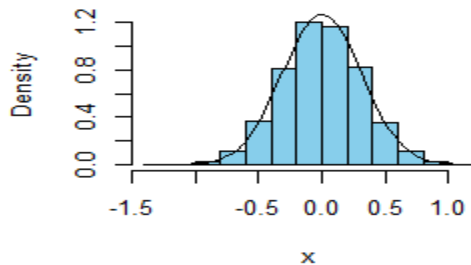




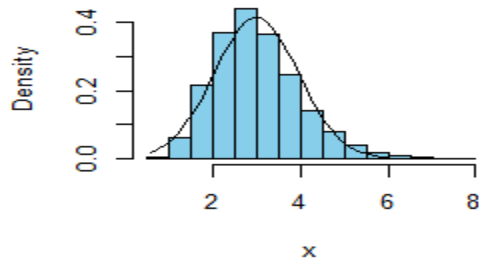
# clt.examp(10)

sample size = 10

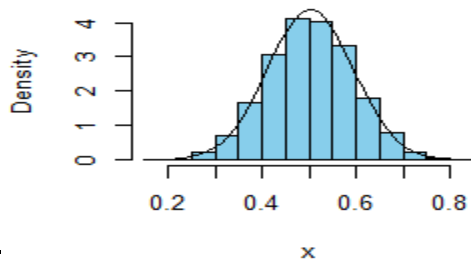
**Normal**



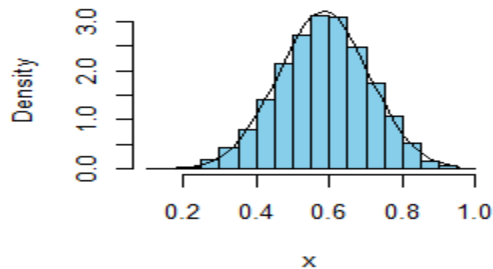
**Gamma**



**Uniform**

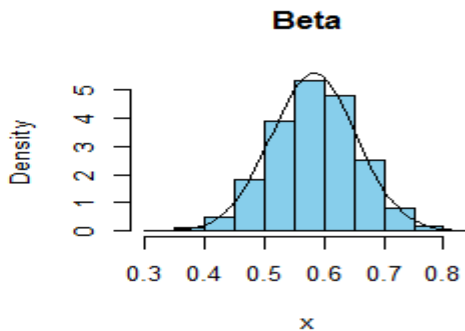
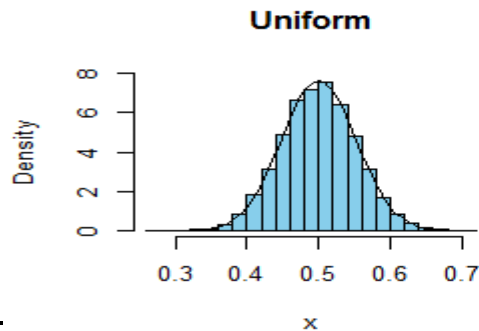
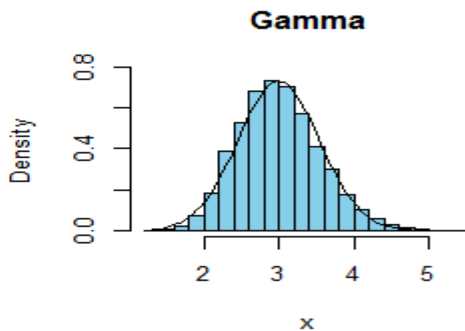
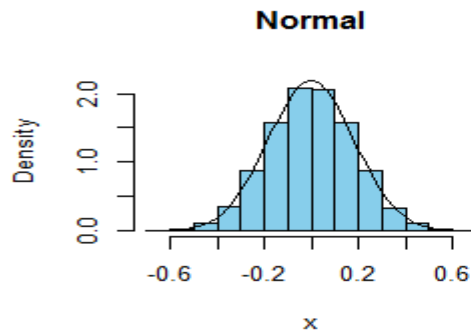


**Beta**



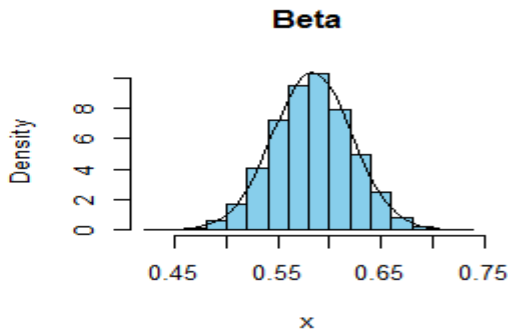
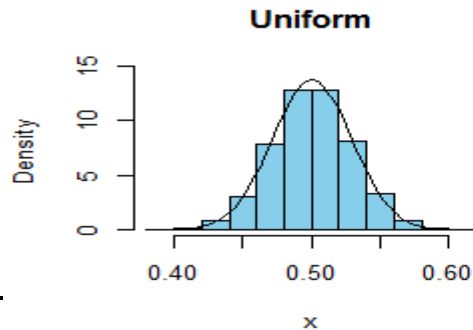
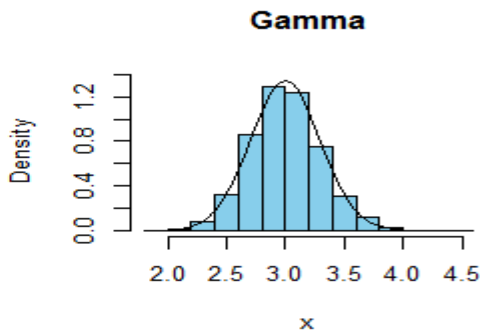
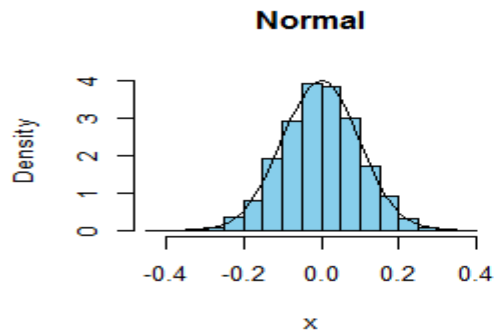
# clt.examp(30)

sample size = 30



# clt.examp(100)

sample size = 100



---

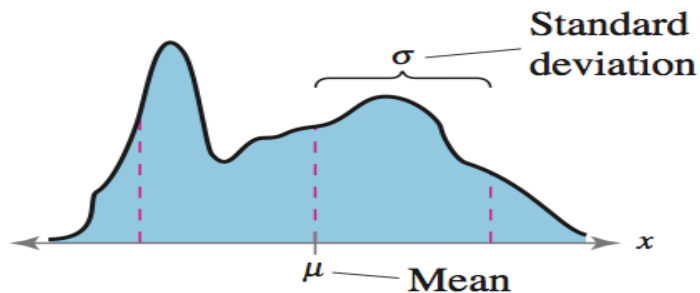
# Code On HW

## ■ We have you examine

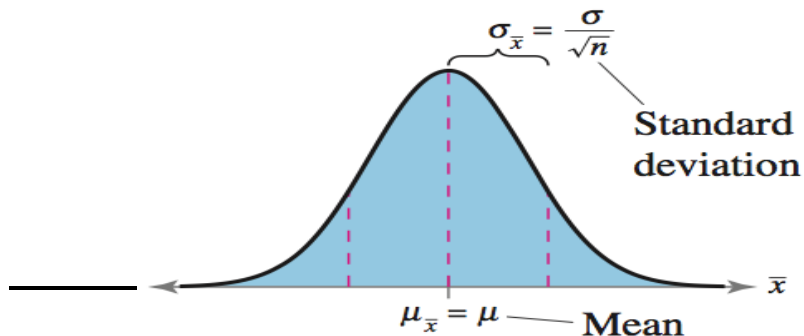
```
simdata=matrix(nrow=10000,ncol=100,runif(100*10000,1))  
hist(apply(simdata,1,mean))
```

# Recap-The Central Limit Theorem

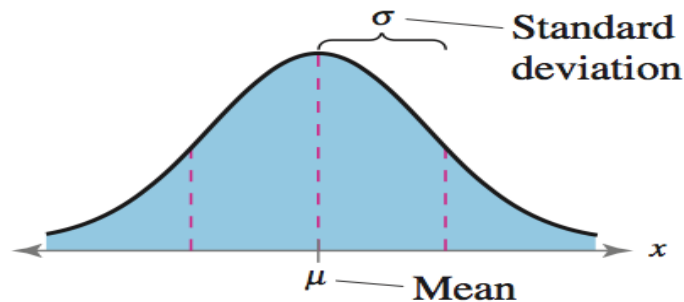
## 1. Any Population Distribution



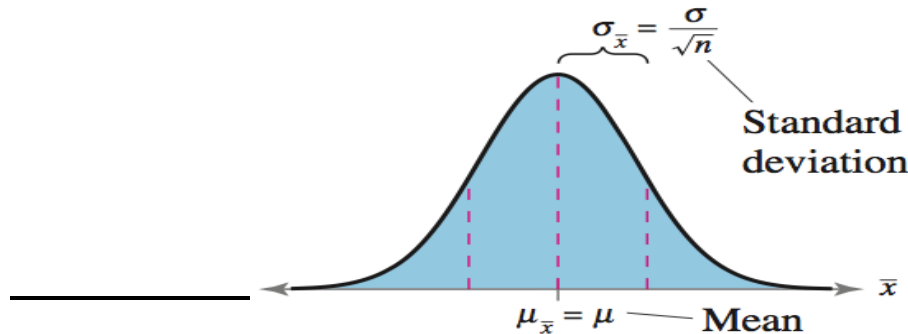
Distribution of Sample Means,  $n \geq 30$



## 2. Normal Population Distribution



Distribution of Sample Means, (any  $n$ )

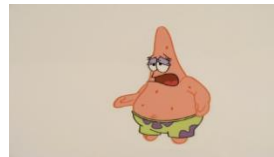


# Example of Using the CLT

- Suppose a population has mean  $\mu = 8$  and standard deviation  $\sigma = 3$ . Suppose a random sample of size  $n = 36$  is selected.
- What is the probability that the sample mean is between 7.8 and 8.2?

# Wait! Lets ask a different question

- What is the probability that a single observation between 7.8 and 8.2?



Do we know if it's a normal distribution?

Do we know if it's a binomial distribution?



With the information we are given, we can only answer questions about the **sample mean** (because of the clt).

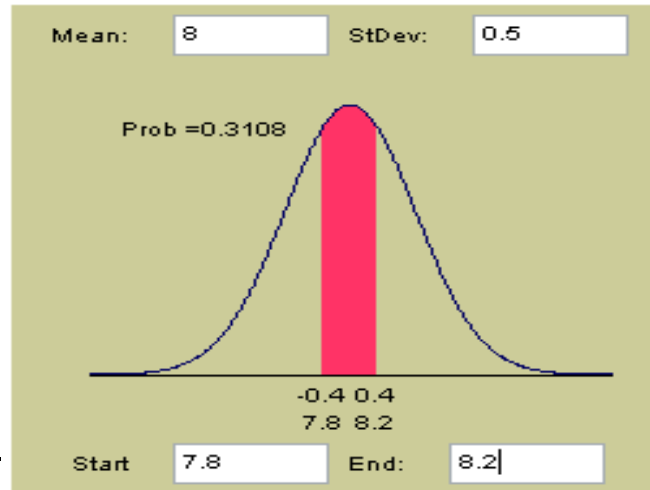
So: What is the probability that the **sample mean** is between 7.8 and 8.2?

- Even if the population is not normally distributed, the central limit theorem can be used ( $n > 30$ )
- ... so the sampling distribution of  $\bar{x}$  is approximately normal
- ... with mean  $\mu_{\bar{x}} = 8$
- ...and standard deviation  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{3}{\sqrt{36}} = 0.5$



# Solution (cont)

$$\begin{aligned} P(7.8 < \bar{X} < 8.2) &= P\left(\frac{7.8-8}{\frac{3}{\sqrt{36}}} < \frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}} < \frac{8.2-8}{\frac{3}{\sqrt{36}}}\right) \\ &= P(-0.4 < Z < 0.4) = 0.3108 \end{aligned}$$



# Recap

- Unless you are explicitly told the distribution of a random variable  $X$ , there is no way to evaluate  $P(a < X < b)$ .
- However, without needing to know the underlying distribution, if  $n$  is sufficiently large, the CLT allows one to evaluate

$$P(a < \bar{X} < b)$$

- This is a powerful result.

# Example of using the CLT

## ■ From HW

Annual real returns on the Standard & Poor's 500-Stock Index over the period 1871 to 2004 have varied with mean 9.2% and standard deviation 20.6%. Andrew plans to retire in 45 years and is considering investing in stocks.

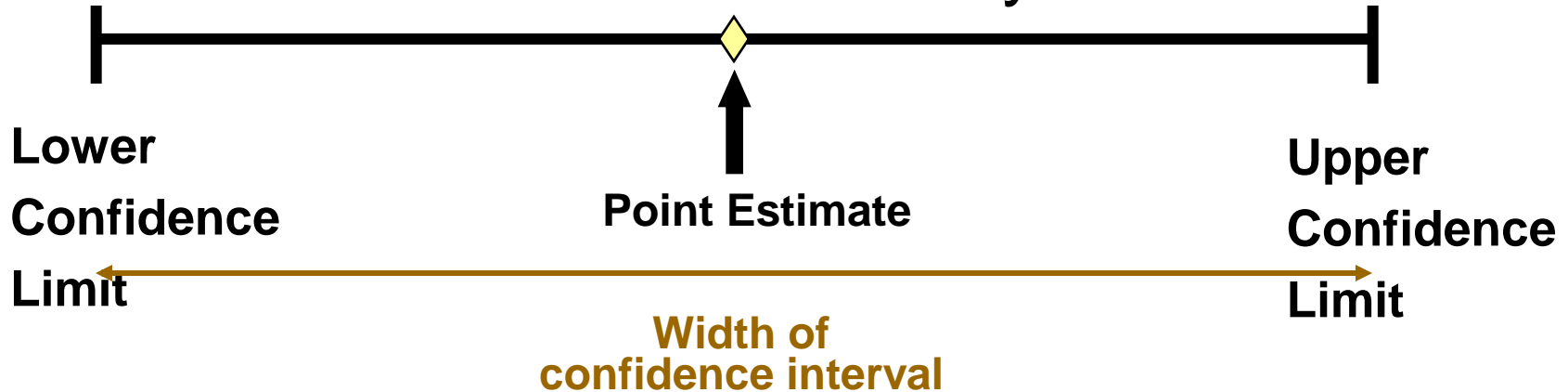
- a) What is the probability (assuming that the past pattern of variation continues) that the mean annual return on common stocks over the next 45 years will exceed 15%?
- b) What is the probability that the mean return will be less than 5%?

# How Large is Large Enough?

- For most distributions,  $n > 30$  will give a sampling distribution that is nearly normal
- For fairly symmetric distributions,  $n > 15$
- For normal population distributions, the sampling distribution of the mean is always normally distributed

# Confidence Intervals

- A **point estimate** is a single number,
- a **confidence interval** provides additional information about variability



# CI for the Mean

- Derived from the Central Limit Theorem
- Assuming  $n > 30$ , a 95% confidence interval for the population mean  $\mu$  is given by

$$\bar{x} \pm 1.96 \frac{s}{\sqrt{n}}$$

# Example

- Suppose over the last four years the mean return of FBALX is 0.0078 with a sd of 0.025.
- We are then 95% the true mean return is somewhere in the interval

```
> c(.0078-1.96*.025/sqrt(48), .0078+1.96*.025/sqrt(48))  
[1] 0.0007274592 0.0148725408
```

# CI for a Proportion

- Derived from the Central Limit Theorem
- Assuming  $n > 30$ , a 95% confidence interval for the population proportion  $p$  is given by

$$\hat{p} \pm 1.96 \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$



# Example

- Over the last 48 months, AAPL has had 26 positive returns. Give a 95% confidence interval for the true proportion of monthly positive returns.

```
> phat=26/48  
> c(phat-1.96*sqrt(phat*(1-phat)/48),phat+1.96*sqrt(phat*(1-phat)/48))  
[1] 0.4007079 0.6826255
```

# Two Sample Confidence Intervals

## ■ For two means

$$(\bar{x}_1 - \bar{x}_2) \pm 1.96 \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

## ■ For two proportions

$$(\hat{p}_1 - \hat{p}_2) \pm 1.96 \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

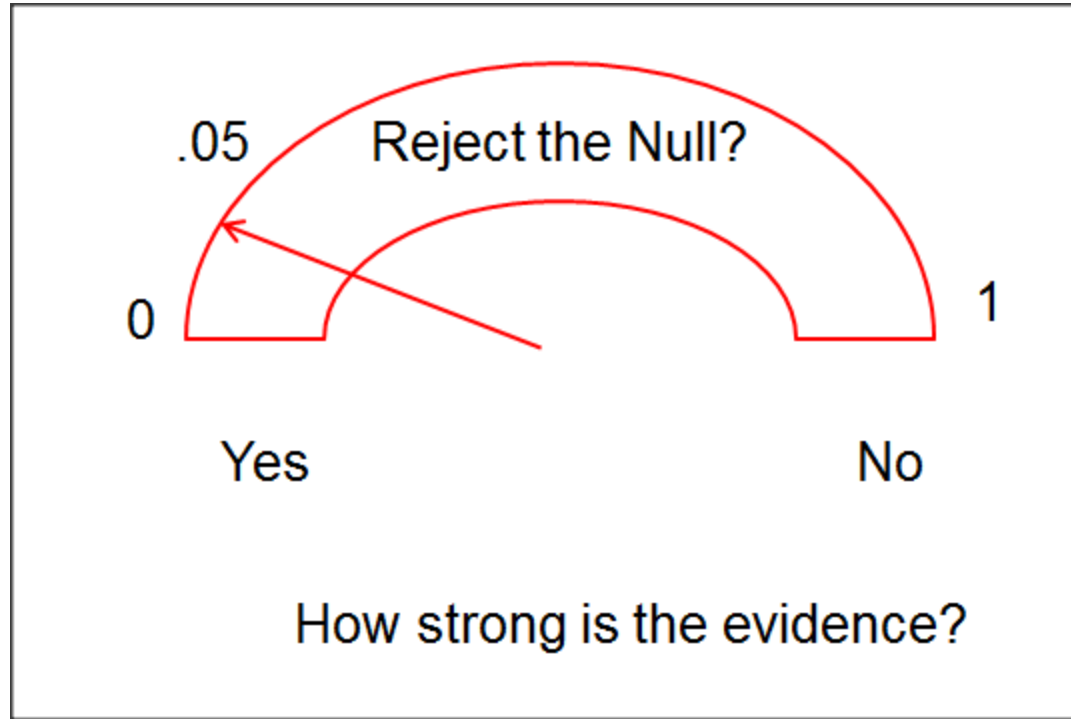
# Hypothesis Tests

- Test some assertion about the population or its parameters
- Can never determine truth or falsity for sure – only get evidence that points one way or another
- *Null hypothesis* ( $H_0$ ) – what is to be tested
- *Alternate hypothesis* ( $H_a$ ) – denial of  $H_0$
- Develop a decision rule to decide on  $H_0$  or  $H_a$  based on sample data

# $p$ -Values for Hypothesis Tests

- $p$ -value quantifies confidence about the decision
  - If  $p < 0.05$ , reject  $H_0$
  - If  $p \geq 0.05$ , do not reject  $H_0$
- The  $p$ -value quantifies how strong the evidence is in favor of the null or alternative hypothesis.
- **Mantra-if  $p$  is low  $H_0$  must go!**

# P-values for Hyp Testing



The  $P$  value measures the strength of evidence against the null hypothesis.

# Hypothesis Testing (in one slide!!!)

- The usual hypothesis testing set-up is as follows:

$H_0 : A \text{ is true}$

The null  
hypothesis

$H_a : B \text{ is true}$

The alternative  
hypothesis

- We reject the null hypothesis if the  
p-value < .05.

# Example: Testing Normality

- There are many ways to test if a data set comes from a normal distribution
- `shapiro.test`
- `jaque.bera.test` (in package `tseries`)
- `ad.test` (in package `normtest`)

$H_0$  : data is normal

$H_a$  : data is not normal

# Normality Test

- If the tails are fat, the distribution isn't normal. This can be tested using a normality test.
- $H_0$ : data normal    $H_a$ : data not normal

```
> shapiro.test(rnorm(100))
```

```
Shapiro-Wilk normality test
```

```
data:  rnorm(100)  
W = 0.99004, p-value = 0.6678
```

- If p is low  $H_0$  must go



# Normality Test note as.numeric()

## ■ Interesting, time period matters

```
> shapiro.test(as.numeric(monthlyReturn(Ad(JPM))))
```

Shapiro-Wilk normality test

```
data:  as.numeric(monthlyReturn(Ad(JPM)))  
W = 0.98867, p-value = 0.5227
```

```
> shapiro.test(as.numeric(dailyReturn(Ad(JPM))))
```

Shapiro-Wilk normality test

```
data:  as.numeric(dailyReturn(Ad(JPM)))  
W = 0.82996, p-value < 2.2e-16
```

# Testing Means and Proportions

- On the next few pages we quickly work through some examples.
- More examples are found under hw2 on the website.

# Testing Means

■ Do FCNTX and FDVLX have the same average return?

```
> t.test(monthlyReturn(FDVLX), monthlyReturn(FCNTX))
```

Welch Two Sample t-test

```
data:  monthlyReturn(FDVLX) and monthlyReturn(FCNTX)
t = -0.07129, df = 207.42, p-value = 0.9432
alternative hypothesis: true difference in means is not
equal to 0
95 percent confidence interval:
 -0.01427305  0.01327682
sample estimates:
 mean of x    mean of y 
0.004201586 0.004699701
```

# Testing Variances

## ■ Do FCNTX and FDVLX have the same variance?

```
> var.test(monthlyReturn(FDVLX), monthlyReturn(FCNTX))
```

```
      F test to compare two variances
```

```
data:  monthlyReturn(FDVLX) and monthlyReturn(FCNTX)
F = 2.0498, num df = 116, denom df = 116, p-value = 0.0001352
alternative hypothesis: true ratio of variances is not equal to 1
95 percent confidence interval:
 1.421927 2.955047
sample estimates:
              monthly.returns
monthly.returns      2.049844
attr(,"names")
[1] "ratio of variances"
```

# Testing Correlations

```
> cor.test(monthlyReturn(FDVLX), monthlyReturn(FCNTX))
```

Pearson's product-moment correlation

```
data:  monthlyReturn(FDVLX) and monthlyReturn(FCNTX)
t = 19.747, df = 115, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to
0
95 percent confidence interval:
 0.8295836 0.9144323
sample estimates:
      cor
0.8787758
```

# Density Estimation..Why be Normal?

- Density Estimation refers to determining a probability density for a given set of data.
- This can be as simple as assuming the data is normal, and using the sample mean and standard deviation as parameters.
- Or it can involve some fancy math.

# What do densities require?

- A probability density is a continuous function  $f(x)$  so that
- $f(x) \geq 0$  for all values of  $x$
- The function  $f(x)$  integrates to 1.

# The logspline package (unique to R)

- Library(logspline)
- Performs a semi-parametric density estimate, then lets you sample from the estimate.
- Pretty cool routine.
- What????



# The Stone-Weierstrass Theorem

## ■ From Wiki:

In mathematical analysis, the **Weierstrass approximation theorem** states that every continuous function defined on an interval  $[a,b]$  can be uniformly approximated as closely as desired by a polynomial function. Because polynomials are among the simplest functions, and because computers can directly evaluate polynomials, this theorem has both practical and theoretical relevance, especially in polynomial interpolation. The original version of this result was established by Karl Weierstrass in 1885.

# Example

- Consider the function  $\sin(x)$
- Consider also the function  $P(x)$

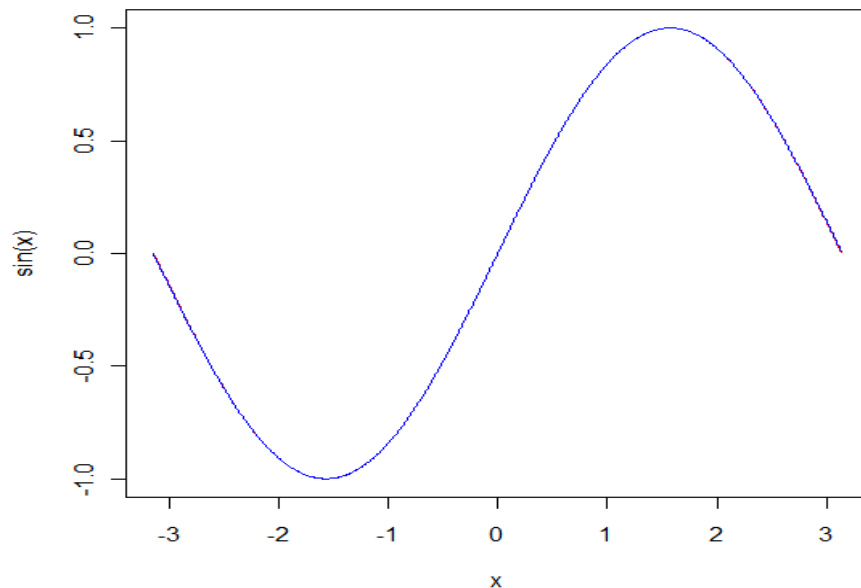
$$P(x) = x - \frac{1}{6}x^3 + \frac{1}{120}x^5 - \frac{1}{5040}x^7 + \frac{1}{362,880}x^9$$

- On the next slide is a graph of these two functions.

# Can't tell them apart

```
> plot(x, sin(x), col="red", type="l")
```

```
> lines(x, x-x^3/6+x^5/120-x^7/5040+x^9/362880, col="blue")
```



# Approximate Densities

- We can do a similar approximation for any density.
- That is, given a set of data, we can find a polynomial approximation for the density of the underlying distribution.
- This is a powerful technique that we will use when we do Monte Carlo simulation.

# The logspline package

- The package in R is called logspline
- It allows one to fit a density estimate then
  - Calculate Probabilities
  - Generate Random Observations
- Very powerful routine
  - `logspline()` – fit the density estimate
  - `dlogspline()` - density
  - `rlogspline()` – generate random values

# Example

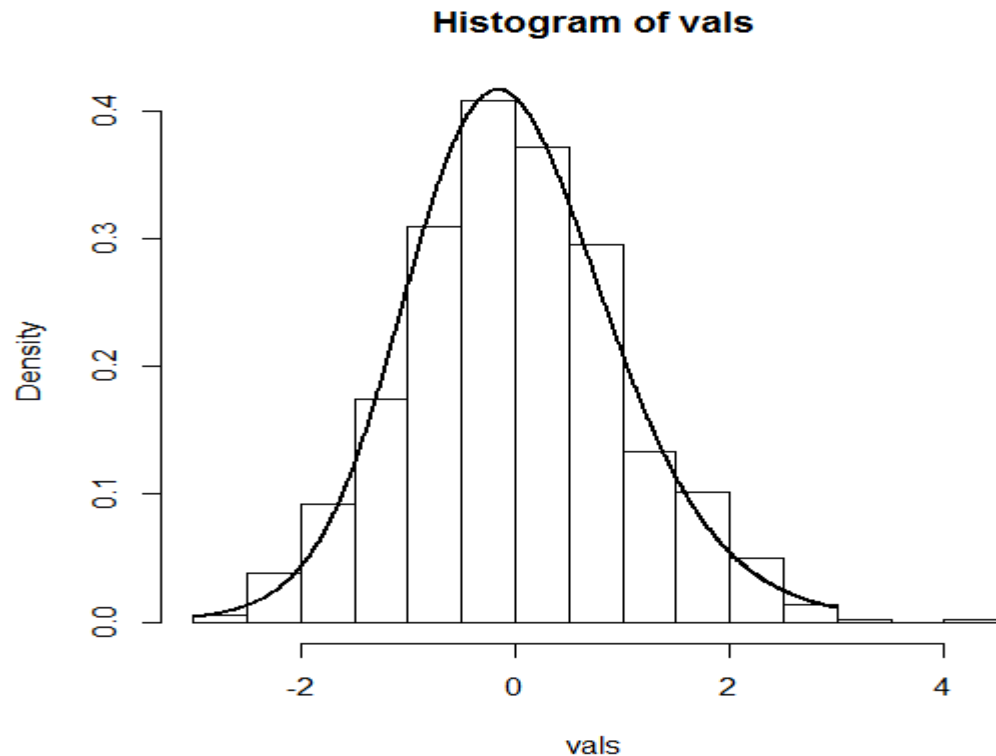
## ■ Fit a density to some data

```
> vals=rnorm(100)
> fit=logspline(vals)
```

## ■ Draw a histogram and overlay the estimated density.

```
> hist(vals,prob=TRUE)
> x=seq(-3,3,.01)
> lines(x,dlogspline(x,fit))
```

# The Estimated Density



From a histogram of 100 values, it comes close to reproducing the underlying density

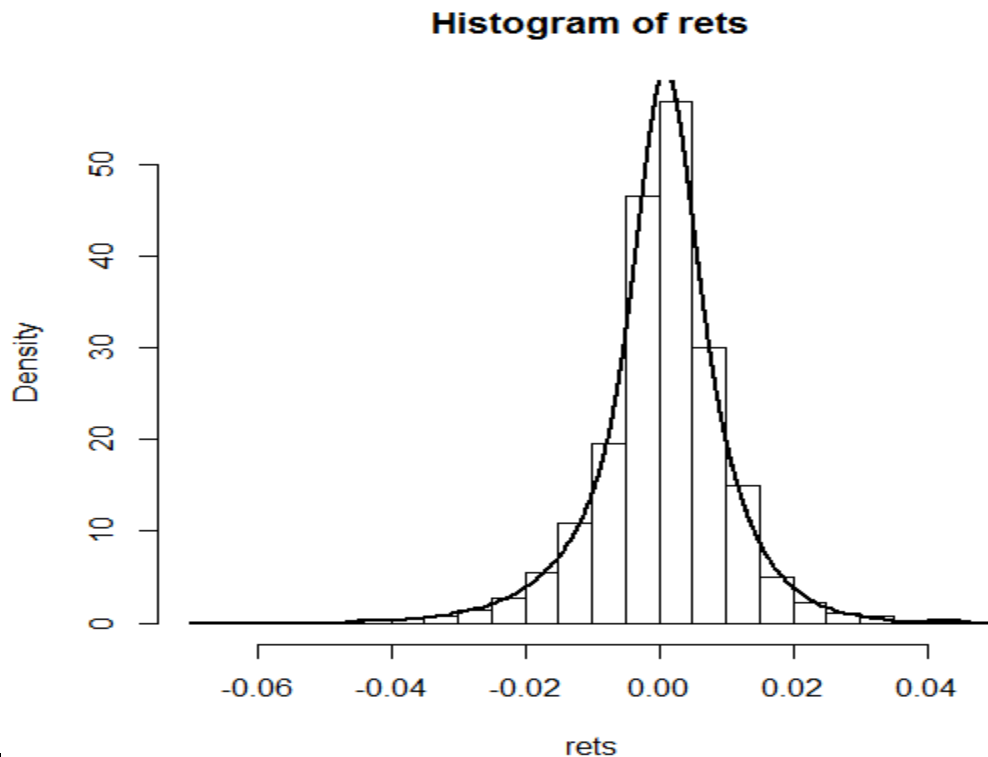
# Model historical stock returns

## ■ Just a few lines of code

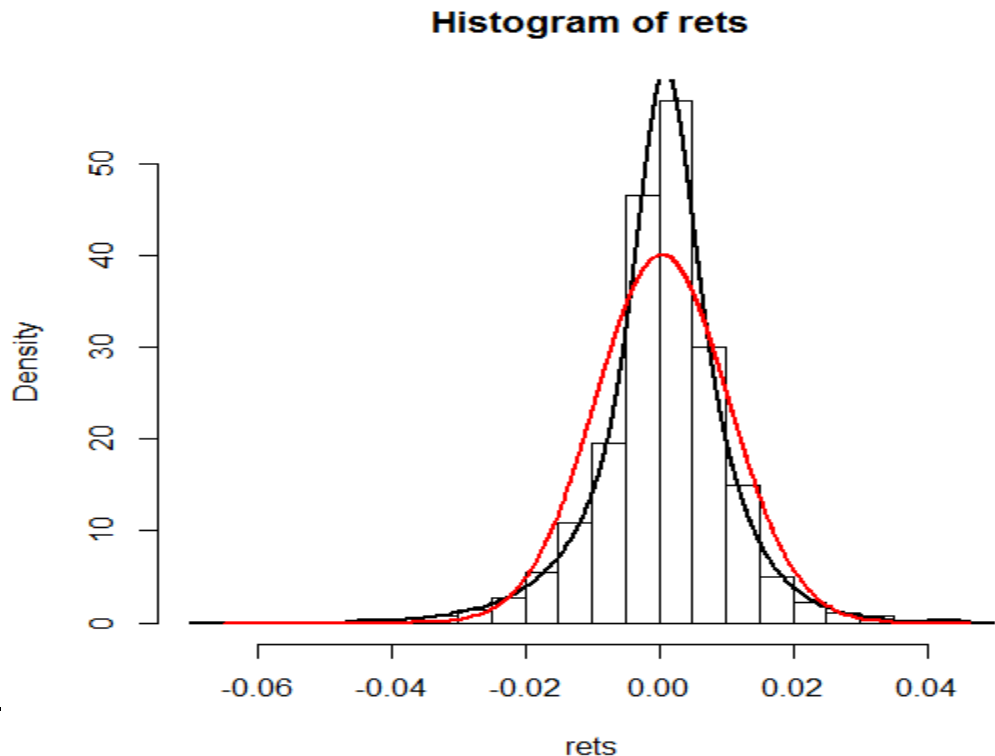
```
> getSymbols("SPY", from="2010-01-01")  
[1] "SPY"  
> rets=dailyReturn(Ad(SPY))  
> rets=as.numeric(rets)  
> fit=logspline(rets)
```



# The resulting density



# Compare with the Normal



**Don't make  
assumptions like  
normality! Let your data  
do the talking.**

# Compare Probabilities

- What is the probability of a negative daily return?

```
> pnorm(0, mean(rets), sd(rets))  
[1] 0.4794895
```

```
> plog spline(0, fit)  
[1] 0.4441395
```

# Can draw from this density

- A powerful ability of the logspline routine is that one **can draw random values** from the fitted density function.
- This is easily done using the rlogspline routine.
- Good way to simulate historical stock returns.
- More details when we cover Monte Carlo Simulation.