

INFO 7374 Machine Learning in Finance

Lecture 1 Basic Statistics.....	3
Probability	3
Common Distributions	7
Uniform.....	7
Normal	8
Log-Normal Distribution	9
Students–t Distribution	11
Exponential	12
Gamma.....	14
Chi-Square	15
Weibull	16
Beta	16
Multivariate Normal	17
Generating Random Variables.....	19
Uniform Random Numbers.....	19
Normal Random Variables	19
Jointly Normal Random Variables.....	20
Inverse Transform Method	20
Acceptance-Rejection Method.....	21
Metropolis-Hastings Algorithms.....	22
Non-parametric Kernel Density	23
Kernel Density Estimator.....	23
Kernel Smoothing Function	24
Bootstrap.....	27
Lecture 2 Machine Learning in Market Making	29
Market Making	29
Algorithm Trading	29
Choice of Quantitative Models.....	31
Parameter Estimation Methods.....	31
Construction of Test Statistics	31
Review on Time-Series Econometrics.....	31
Definitions for Nonstationary Time Series	32
Unit Root	32
Cointegration	35
Paris Trading.....	35
Conditional Heteroskedasticity	36
Filtering Methods	37
The Kalman Filter.....	38
Hamilton Switching Model.....	39
Particle Filtering	40
Support Vector Machine.....	41
Linear SVM Regression: Dual Formula	43
Nonlinear SVM Regression: Primal Formula	44
Predictive Accuracy.....	47
Diebold Mariano Test	47
Amisano Giacomini Test	48
Back Test Trading Strategies.....	48
Cross Validation	48

Maximum Likelihood Estimation	49
Lecture 3 Machine Learning in Active Investment	50
A Short History of Investment Beliefs	50
Quantitative Methods Involved	50
Active Alpha Strategies	50
Regression Analysis	52
Ordinary Least Square Estimation	52
When Matrices Are Singular	55
Pitfalls with Regression	55
Multicollinearity	55
Heteroscedasticity	58
Autocorrelation	63
Model Selection	67
Akaike Information Criterion	67
Bayesian Information Criterion	67
Ridge Regression	67
LASSO Regression	68
Least Angle Regression (LARS).....	68
Elastic Net.....	70
Principle Components Analysis	70
Examples of Machine Learning in Active Investment.....	71
Neural Networks	77
Random Forecast	78
Stochastic Gradient Descent	79
Lecture 4 Machine Learning in Passive Investment.....	80
Smart Beta Strategy	80
Unsupervised Learning and Smart Beta	81
Hierarchical Clustering	81
k-Means Clustering	82
Model-Based Clustering	84
Lecture 5 Machine Learning in Optimization.....	88
Parameter Estimation.....	88
Portfolio Optimization	89
Unconstrained Optimization	91
Decent Methods	91
Exact Line Search	92
Backtracking Line Search	92
Gradient Descent Method	93
Steepest Descent Method.....	94
Newton's Method	95
Quasi Newton Method	96
Interior-point Methods	96
A Witty Little Technique	98
Genetic Programming	99
Lecture 6 Machine Learning in Risk Management	103
Market Risk Management.....	103
Credit Risk Management	111
Probability of Default	111
Overview of LGD and EAD Modeling Methodologies	111

Lecture 1 Basic Statistics**Probability**

- A **random experiment** is defined as a process or action whose outcome cannot be predicted with certainty and would likely change when the experiment is repeated.
- The **sample space** is the set of all outcomes from an experiment.
- The outcomes from random experiments are often represented by an uppercase variable such as X . This is called a **random variable**, and its value is subject to the uncertainty intrinsic to the experiment.
- Random variables can be discrete or continuous. A **discrete random variable** can take on values from a finite or countably infinite set of numbers. Examples of discrete random variables are the number of defective parts or the number of typographical errors on a page. A **continuous random variable** is one that can take on values from an interval of real numbers.

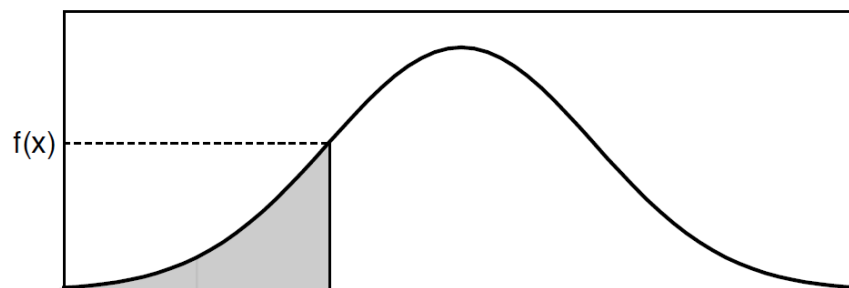
Discrete Random Variables: Letting 1 represent the *bull market* and letting 0 represent the *bear market*, then the probability of the event that *we are in the bull market* would be written as

$$P(X = 1)$$

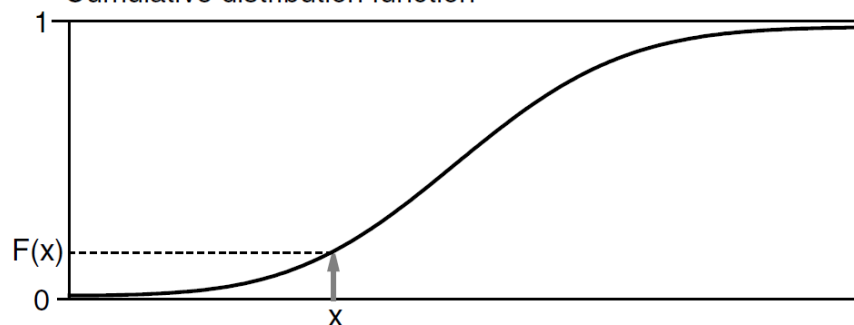
Continuous Random Variables: Let X denote the price of crude oil future: \$/barrel. The probability that the transaction price is in the range \$30 to \$50 is expressed as

$$P(\$30 < X \leq \$50).$$

Probability density function



Cumulative distribution function



- The **mean** or **expected value** of a random variable is defined using the probability density (mass) function. It provides a measure of central tendency of the distribution. If we observe many values of the random variable and take the average of them, we would expect that value to be close to the mean. The expected value is defined below for the discrete case.

EXPECTED VALUE - DISCRETE RANDOM VARIABLES

$$\mu = E[X] = \sum_{i=1}^{\infty} x_i f(x_i)$$

We see from the definition that the expected value is a sum of all possible values of the random variable where each one is weighted by the probability that X will take on that value.

- The **variance** of a discrete random variable is given by the following definition.

VARIANCE - DISCRETE RANDOM VARIABLES for $\mu < \infty$,

$$\sigma^2 = V(X) = E[(X - \mu)^2] = \sum_{i=1}^{\infty} (x_i - \mu)^2 f(x_i).$$

We see that the variance is the sum of the squared distances, each one weighted by the probability that $X = x_i$. Variance is a measure of dispersion in the distribution. If a random variable has a large variance, then an observed value of the random variable is more likely to be far from the mean μ . The standard deviation σ is the square root of the variance.

The mean and variance for continuous random variables are defined similarly, with the summation replaced by an integral. The mean and variance of a continuous random variable are given below.

EXPECTED VALUE - CONTINUOUS RANDOM VARIABLES

$$\mu = E[X] = \int_{-\infty}^{\infty} x f(x) dx.$$

VARIANCE - CONTINUOUS RANDOM VARIABLES for $\mu < \infty$,

$$\sigma^2 = V(X) = E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx.$$

In addition $V(X) = E[X^2] - \mu^2 = E[X^2] - (E[X])^2$.

- Other expected values that are of interest in statistics are the moments of a random variable. These are the expectation of powers of the random variable. In general, we define the r -th moment as

$$\mu'_r = E[X^r]$$

and the **r th central moment** as

$$\mu_r = E[(X - \mu)^r]$$

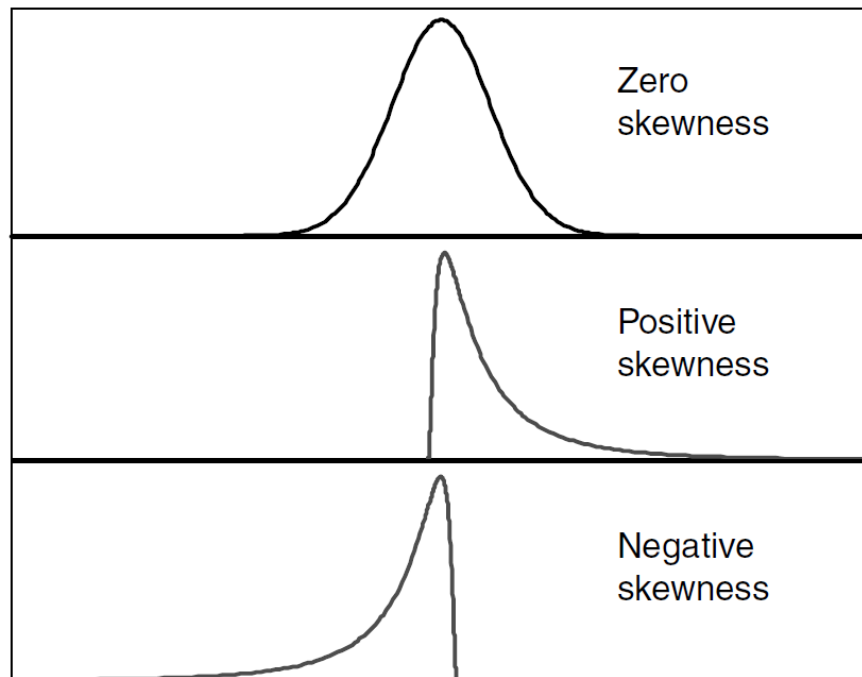
The mean corresponds to μ_1 and the variance is given by μ_2 .

- **Skewness:** The third central moment μ_3 is often called a measure of asymmetry or skewness in the distribution. The uniform and the normal distribution are examples of symmetric distributions. The gamma and the exponential are examples of skewed or asymmetric distributions. The following ratio is called the **coefficient of skewness**, which is often used to measure this characteristic

$$\gamma_1 = \frac{\mu_3}{\mu_2^{3/2}} = \left(\int_{-\infty}^{\infty} (x - E(X))^3 f(x) dx \right) / \sigma^3$$

Distributions that are skewed to the left will have a negative coefficient of skewness, and distributions that are skewed to the right will have a positive value. The coefficient of skewness is zero for symmetric distributions. However, a coefficient of skewness equal to zero does not mean that the distribution must be symmetric..

Probability density function



- **Kurtosis:** Kurtosis measures a different type of departure from normality by indicating the extent of the peak (or the degree of flatness near its center) in a distribution. The **coefficient of kurtosis** is given by the following ratio:

$$\gamma_2 = \frac{\mu_4}{\mu_2^2} = \left(\int_{-\infty}^{\infty} (x - E(X))^4 f(x) dx \right) / \sigma^4$$

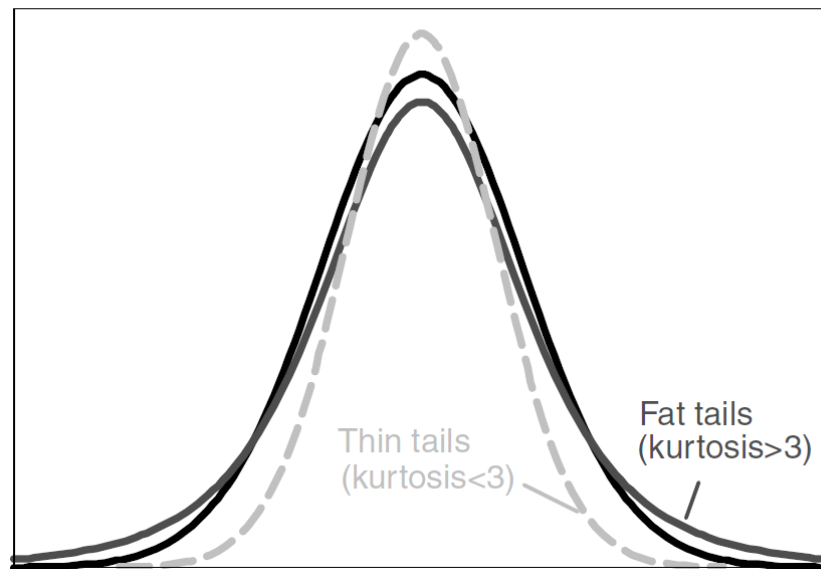
We see that this is the ratio of the fourth central moment divided by the square of the variance. because of the fourth power, large observations in the tail will have a larger weight and hence create large kurtosis. Such a distribution is called **leptokurtic** or **fat-tailed**. If the distribution is normal, then this ratio is equal to 3. A ratio greater than 3 indicates more values in the neighborhood of the mean (is more peaked than the normal distribution). If the ratio is less than 3, then it is an indication that the curve is flatter than the normal.

Sometimes the **coefficient of excess kurtosis** is used as a measure of kurtosis. This is given by

$$\gamma_2 = \frac{\mu_4}{\mu_2^2} - 3.$$

In this case, distributions that are more peaked than the normal correspond to a positive value of γ_2 , and those with a flatter top have a negative coefficient of excess kurtosis.

Probability density function



- The distribution can also be described by its **Quantile**, which is the cutoff point x with an associated probability c :

$$F(x) = \int_{-\infty}^x f(u) du = c$$

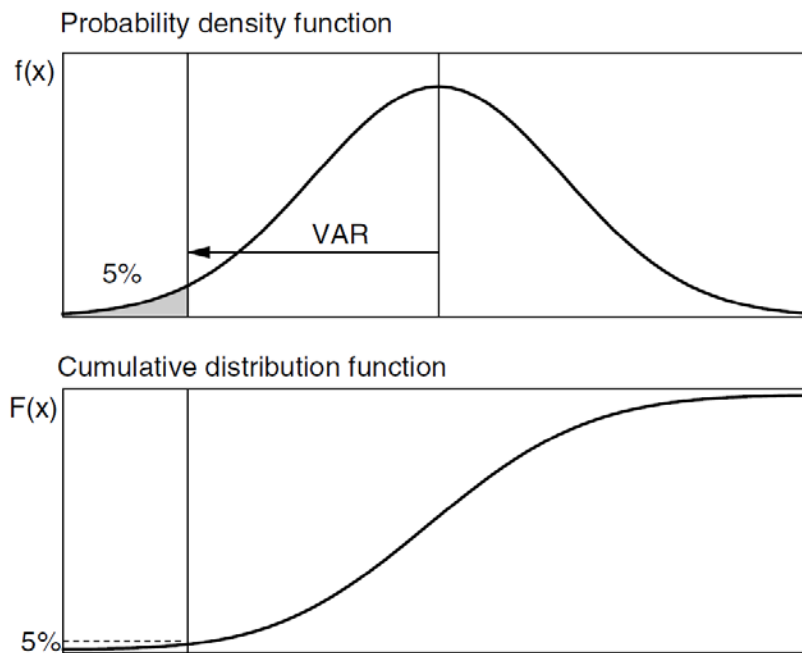
So, there is a probability of c that the random variable will fall below x . Because the total probability adds up to one, there is a probability of $p = 1 - c$ that the random variable will fall above x . Define this quantile as $Q(X, c)$. The 50% quantile is known as the **median**.

Value At Risk (VAR) can be interpreted as the cutoff point such that a loss will not happen with probability greater than $p = 95\%$ percent, say. If $f(u)$ is the distribution of profit and losses on the portfolio, VAR is defined from

$$F(x) = \int_{-\infty}^x f(u)du = 1 - p$$

where p is the right-tail probability, and c usual left-tail probability. VAR can then be defined as the deviation between the expected value and the quantile,

$$VAR(c) = E(X) - Q(X, c)$$



Common Distributions

Uniform

Perhaps one of the most important distributions is the uniform distribution for continuous random variables. One reason is that the uniform $(0,1)$ distribution is used as the basis for simulating most random variables.

A random variable that is uniformly distributed over the interval (a, b) follows the probability density function given by

$$f(x; a, b) = \frac{1}{b - a}$$

$$a < x < b$$

The parameters for the uniform are the interval endpoints, a and b . The mean and variance of a uniform random variable are given by

$$E[X] = \frac{a + b}{2}$$

and

$$V(X) = \frac{(b-a)^2}{12}$$

The cumulative distribution function for a uniform random variable is

$$F(x) = \begin{cases} 0 & x \leq a \\ \frac{x-a}{b-a} & a < x < b \\ 1 & x \geq b \end{cases}$$

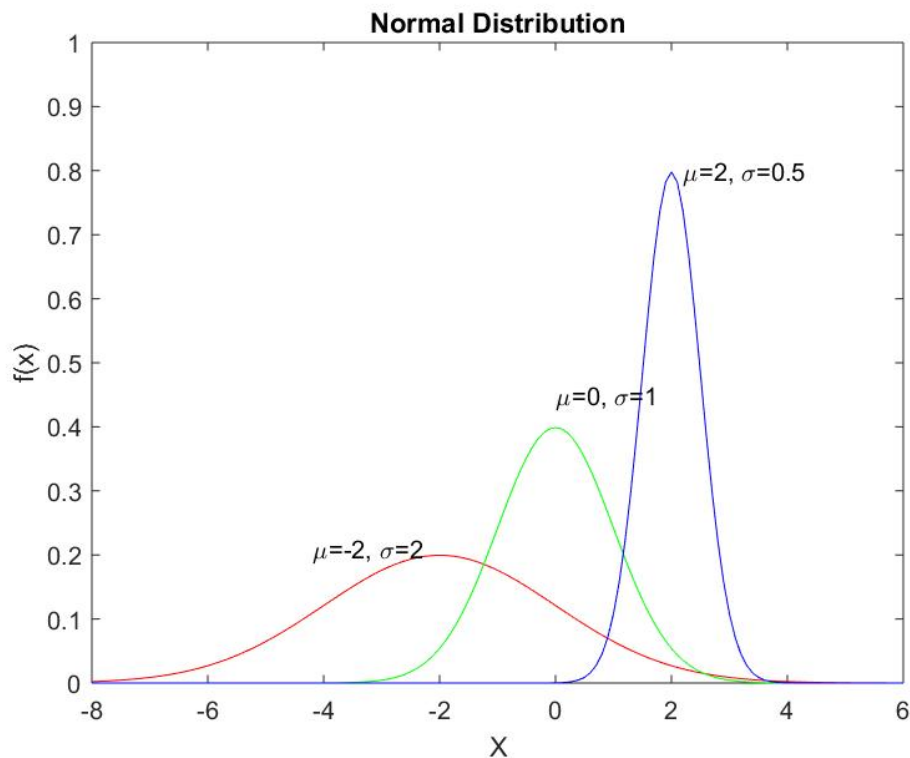
Normal

A well known distribution in statistics and engineering is the normal distribution. Also called the Gaussian distribution, it has a continuous probability density function given by

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

where The normal distribution is completely determined by its parameters (μ and σ^2), which are also the expected value and variance for a normal random variable. The notation $X \sim N(\mu, \sigma^2)$ is used to indicate that a random variable X is normally distributed with mean μ and variance σ^2 .

Several normal distributions with different parameters are shown below.



Some special properties of the normal distribution are given here.

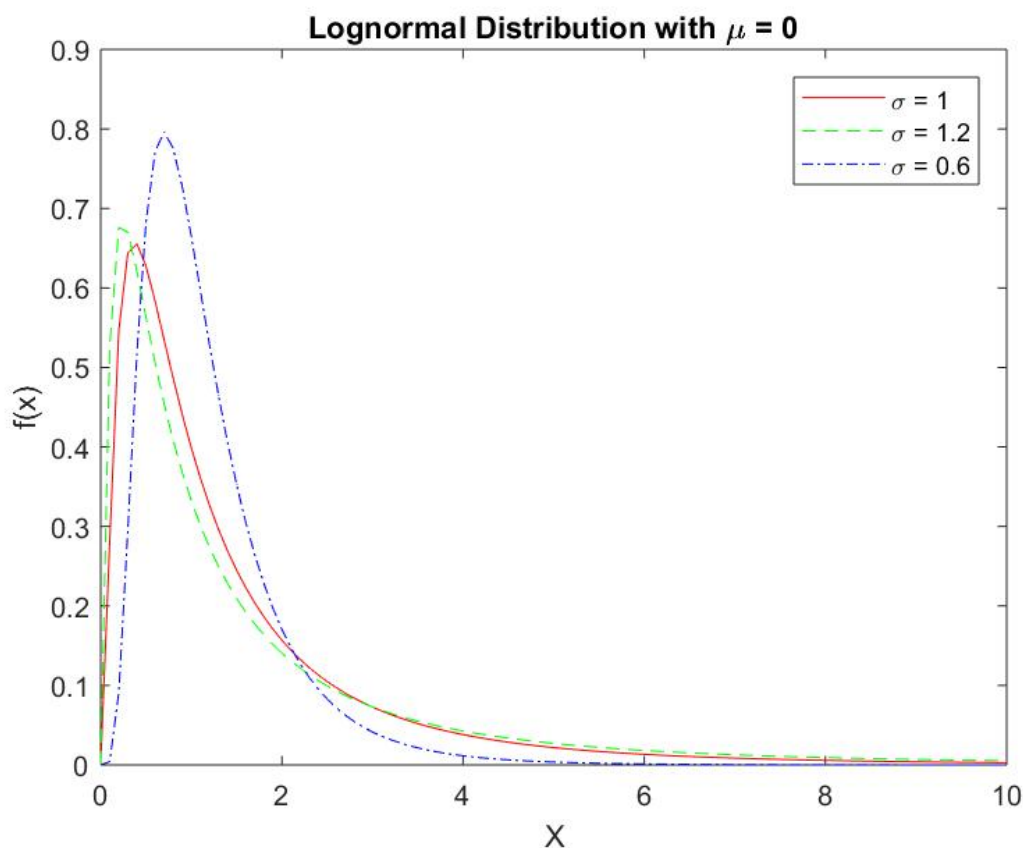
- The value of the probability density function approaches zero as x approaches positive and negative infinity.
- The probability density function is centered at the mean, and the maximum value of the function occurs at $x = \mu$.
- The probability density function for the normal distribution is symmetric about the mean μ .

The special case of a standard normal random variable is one whose mean is zero ($\mu = 0$), and whose standard deviation is one ($\sigma = 1$). If X is normally distributed, then

$$Z = \frac{X - \mu}{\sigma} \sim N(0,1)$$

is a standard normal random variable.

Log-Normal Distribution



The normal distribution is a good approximation for many financial variables, such as the rate of return on a stock, $r = (P_1 - P_0)/P_0$, where P_0 and P_1 are the stock prices at time 0 and 1.

Strictly speaking, this is inconsistent with reality since a normal variable has infinite tails on both sides. Due to the limited liability of corporations, stock prices cannot turn negative. This rules out returns lower than minus unity and distributions with infinite left tails, such as the normal distribution. In many situations, however, this is an excellent approximation. For instance, with short horizons or small price moves, the probability of having a negative price is so small as to be negligible.

If this is not the case, we need to resort to other distributions that prevent prices from going negative. One such distribution is the lognormal.

A random variable X is said to have a **lognormal distribution** if its logarithm $Y = \ln(X)$ is **normally distributed**. This is often used for continuously compounded returns, defining $Y = \ln(P_1/P_0)$. Because the argument X in the logarithm function must be positive, the price can never go below zero. Large and negative large values of Y correspond to P_1 converging to, but staying above, zero.

The lognormal density function has the following expression

$$f(x) = f(x; \mu, \sigma^2) = \frac{1}{x\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(\ln(x) - \mu)^2}{2\sigma^2}\right\}, \quad x > 0$$

Note that this is more complex than simply plugging $\ln(x)$ in Equation above, because x also appears in the denominator. Its mean is

$$E[X] = \exp\left[\mu + \frac{1}{2}\sigma^2\right]$$

and variance

$$V[X] = \exp[2\mu + 2\sigma^2] - \exp[2\mu + \sigma^2].$$

The parameters were chosen to correspond to those of the normal variable,

$$E[Y] = E[\ln(X)] = \mu$$

and

$$V[Y] = V[\ln(X)] = \sigma^2.$$

Conversely, if we set $E[X] = \exp[r]$, the mean of the associated normal variable is

$$E[Y] = E[\ln(X)] = (r - \sigma^2/2).$$

This adjustment is also used in the Black-Scholes option valuation model, where the formula involves a trend in $(r - \sigma^2/2)$ for the log-price ratio.

We also note that the **distribution of the bond price, resembles a lognormal distribution**. Using **continuous compounding instead of annual compounding**, the price function is

$$V = 100 \exp(-rT)$$

which implies $\ln\left(\frac{V}{100}\right) = -rT$. Thus if r is normally distributed, V has a lognormal distribution.

Students- t Distribution

Another important distribution is the Student's t -distribution. This arises in hypothesis testing, because it describes the distribution of the ratio of the estimated coefficient to its standard error.

The distribution is characterized by a parameter k known as the degrees of freedom. Its density is

$$f(x) = \frac{\Gamma[(k+1)/2]}{\Gamma(k/2)} \frac{1}{\sqrt{k\pi}} \frac{1}{\left(1 + \frac{x^2}{k}\right)^{(k+1)/2}}$$

where Γ is the gamma function, defined as

$$\Gamma(k) = \int_0^\infty x^{k-1} e^{-x} dx.$$

As k increases, this function converges to the normal p.d.f.

This distribution is symmetrical with mean zero and variance

$$V[X] = \frac{k}{k-2}$$

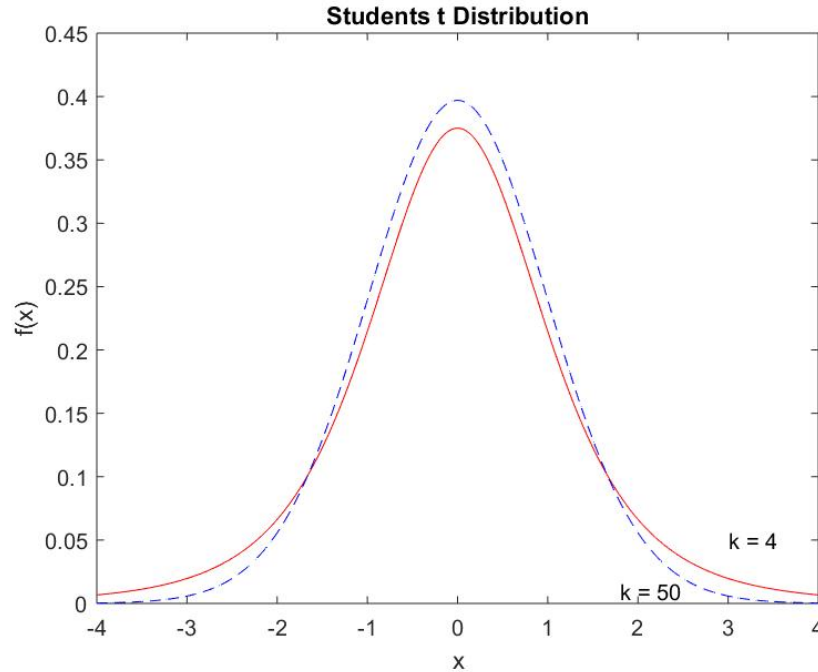
provided $k > 2$. Its kurtosis is

$$\delta = 3 + \frac{6}{k-4}$$

provided $k > 4$. It has fatter tail than the normal distribution, which often provides a better representation of typical financial variables. Typical estimated values of k are around four to six for stock returns. We can also use the Student's t to compute Value-At-Risk as function of volatility.

$$VAR = \alpha_k \sigma$$

where the multiplier now depends on the degrees of freedom k .



Exponential

The **exponential distribution** can be used to model the amount of time until a specific event occurs or to model the time between independent events. Some examples where an exponential distribution could be used as the model are:

- the time until the computer locks up,
- the time between arrivals of telephone calls, or
- the time until a part fails.

The exponential probability density function with parameter λ is

$$f(x; \lambda) = \lambda e^{-\lambda x}$$

$$x \geq 0 ;$$

$$\lambda > 0.$$

The mean and variance of an exponential random variable are given by the following:

$$E[X] = \frac{1}{\lambda},$$

and

$$V(x) = \frac{1}{\lambda^2}.$$

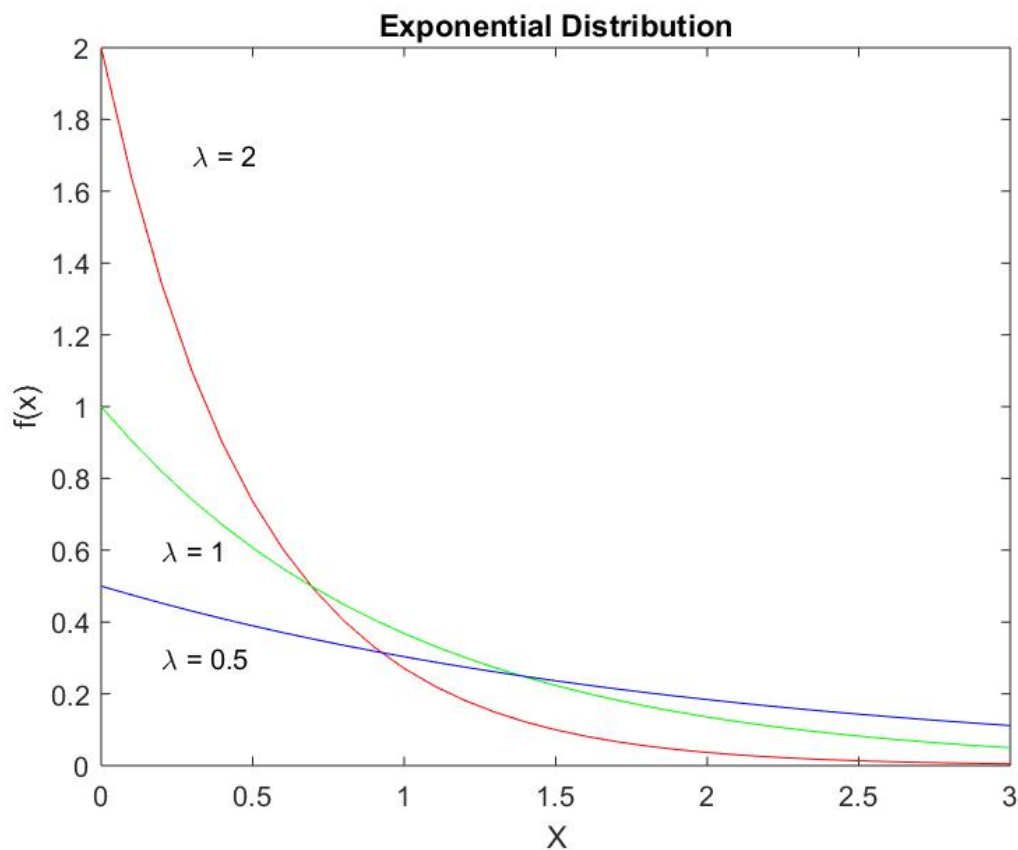
The cumulative distribution function of an exponential random variable is given by

$$F(x) = \begin{cases} 0 & x < 0 \\ 1 - e^{-\lambda x} & x \geq 0 \end{cases}$$

The exponential distribution is the only continuous distribution that has the memory less property. This property describes the fact that the remaining lifetime of an object (whose lifetime follows an exponential distribution) does not depend on the amount of time it has already lived. This property is represented by the following equality, where $s \geq 0$ and $t \geq 0$:

$$P(X > s + t | X > s) = P(X > t)$$

In words, this means that the probability that the object will operate for time, given it has already operated for time s , is simply the probability that it operates for time t . When the exponential is used to represent inter-arrival times, then the parameter is a rate with units of arrivals per time period. When the exponential is used to model the time until a failure occurs, then is the failure rate. Several examples of the exponential distribution are shown in Example below.



Gamma

The gamma probability density function with parameters $\lambda > 0$ and $t > 0$ is

$$f(x; \lambda, t) = \frac{\lambda e^{-\lambda x} (\lambda x)^{t-1}}{\Gamma(t)}$$
$$x \geq 0$$

where t is a shape parameter, and λ is the scale parameter. The gamma function $\Gamma(t)$ is defined as

$$\Gamma(t) = \int_0^{\infty} e^{-y} y^{t-1} dy.$$

For integer values of t , Equation above becomes

$$\Gamma(t) = (t - 1)!.$$

Note that for $t = 1$, the gamma density is the same as the exponential. When t is a positive integer, the gamma distribution can be used to model the amount of time one has to wait until t events have occurred, if the inter-arrival times are exponentially distributed.

The mean and variance of a gamma random variable are

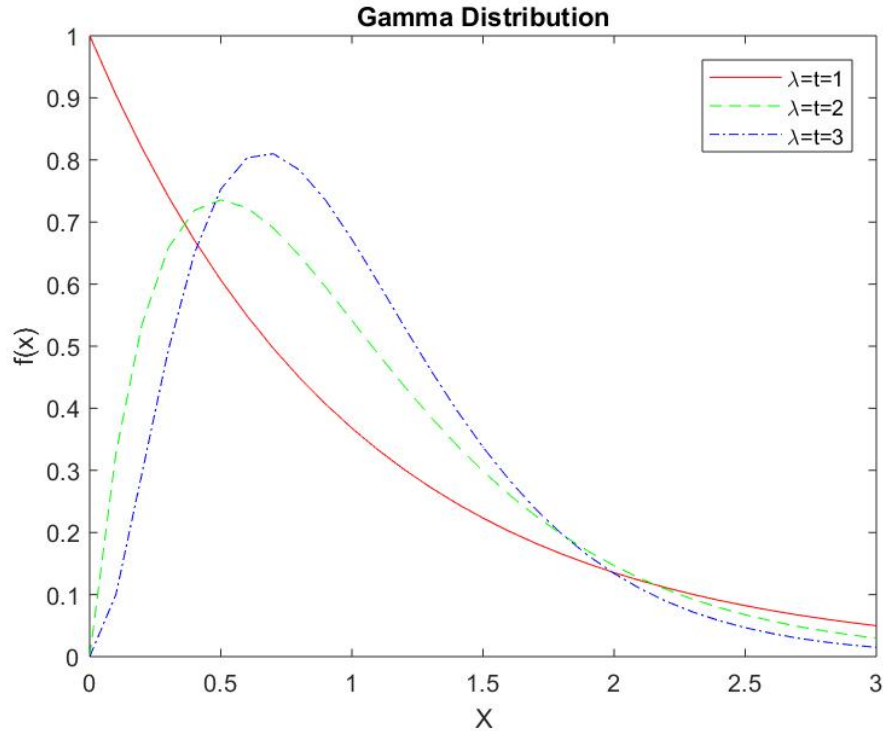
$$E[X] = \frac{t}{\lambda}$$

and

$$V(X) = \frac{t}{\lambda^2}.$$

The cumulative distribution function for a gamma random variable is calculated using

$$F(x; \lambda, t) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{\Gamma(t)} \int_0^{\lambda x} y^{t-1} e^{-y} dy & x > 0 \end{cases}$$



Chi-Square

A gamma distribution where $\lambda = 0.5$ and $t = \frac{\nu}{2}$, with ν a positive integer, is called a chi-square distribution (denoted as χ^2_ν) with ν degrees of freedom. The chi-square distribution is used to derive the distribution of the sample variance and is important for goodness-of-fit tests in statistical analysis.

The probability density function for a chi-square random variable with ν degrees of freedom is

$$f(x; \nu) = \frac{1}{\Gamma(\nu/2)} \left(\frac{1}{2}\right)^{\nu/2} x^{\frac{\nu}{2}-1} e^{-\frac{1}{2}x}$$

$$x \geq 0.$$

The mean and variance of a chi-square random variable can be obtained from the gamma distribution. These are given by

$$E[X] = \nu$$

and

$$V(X) = 2\nu$$

Weibull

The Weibull distribution has many applications in engineering. In particular, it is used in reliability analysis. It can be used to model the distribution of the amount of time it takes for objects to fail. For the special case where $\nu = 0$ and $\beta = 1$, the Weibull reduces to the exponential with $\lambda = 1/\alpha$.

The Weibull density for $\alpha > 0$ and $\beta > 0$ is given by

$$f(x; \nu, \alpha, \beta) = \left(\frac{\beta}{\alpha}\right) \left(\frac{x - \nu}{\alpha}\right)^{\beta-1} e^{-\left(\frac{x-\nu}{\alpha}\right)^\beta} \\ x > \nu$$

and the cumulative distribution is

$$F(x; \nu, \alpha, \beta) = \begin{cases} 0 & x < \nu \\ 1 - e^{-\left(\frac{x-\nu}{\alpha}\right)^\beta} & x \geq \nu \end{cases}$$

The location parameter is denoted by ν and the scale parameter is given by α . The shape of the Weibull distribution is governed by the parameter β . The mean and variance of a random variable from a Weibull distribution are given by

$$E[X] = \nu + \alpha \Gamma\left(\frac{1}{\beta} + 1\right)$$

and

$$V(X) = \alpha^2 \left\{ \Gamma\left(\frac{2}{\beta} + 1\right) - \left[\Gamma\left(\frac{1}{\beta} + 1\right) \right]^2 \right\}$$

Beta

The beta distribution has support on unit interval. It can be used to model a random variable that takes on values over a bounded interval. It has two parameters α and β that determines the shape of the density. A random variable has a beta distribution with parameters $\alpha > 0$ and $\beta > 0$ if its probability density function is given by

$$f(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} \\ 0 < x < 1$$

where

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}.$$

The mean and variance of a beta random variable are

$$E[X] = \frac{\alpha}{\alpha + \beta}$$

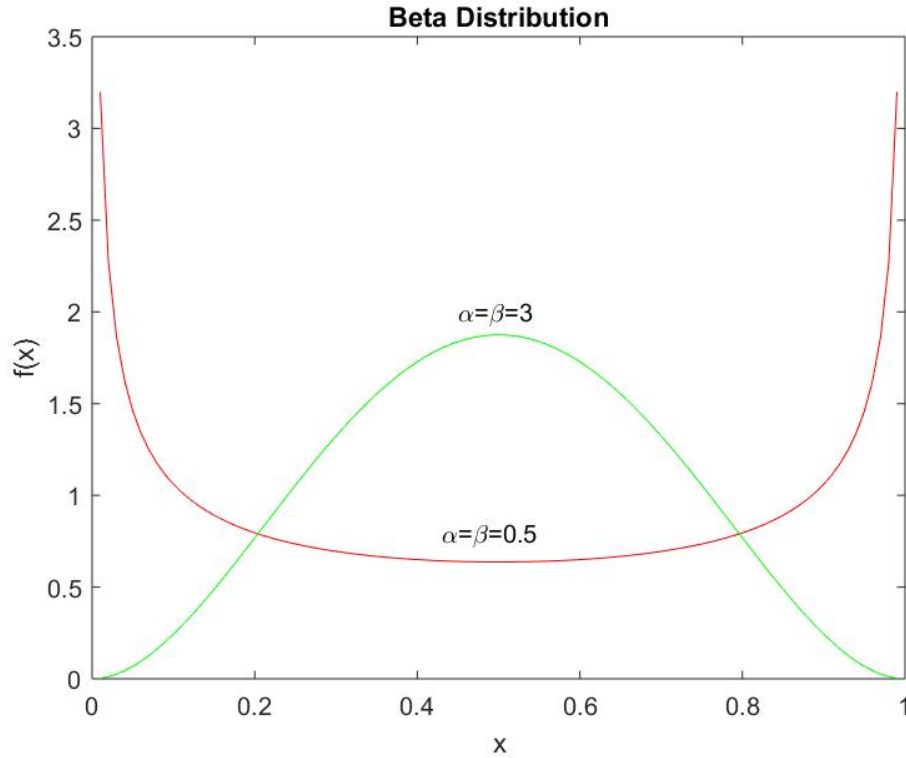
and

$$V(X) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}.$$

The cumulative distribution function for a beta random variable is given by integrating the beta probability density function as follows

$$F(x; \alpha, \beta) = \int_0^x \frac{1}{B(\alpha, \beta)} y^{\alpha} (1 - y)^{\beta-1} dy.$$

The integral in Equation above is called the incomplete beta function.



Multivariate Normal

So far, we have discussed several univariate distributions for discrete and continuous random variables. In this section, we describe one of the important and most commonly used multivariate densities: the multivariate normal

distribution. This distribution is used throughout the rest of the text.

Some examples of where we use it are in exploratory data analysis, in probability density estimation, and in statistical pattern recognition. The probability density function for a general multivariate normal density for d dimensions is given by

$$f(x, \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right\}$$

where x is a d -component column vector, μ is the $d \times 1$ column vector of means, and Σ is the $d \times d$ covariance matrix.

For example, when $d = 2$, $\mu = (\mu_1, \mu_2)^T$ and $\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$

The superscript T represents the transpose of an array.

The mean and covariance are calculated using the following formulas:

$$\mu = E[x],$$

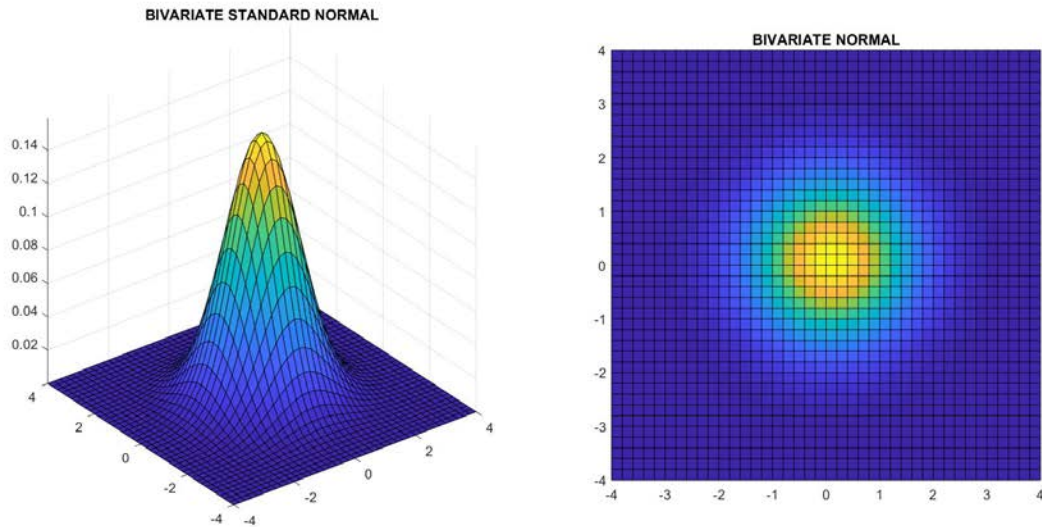
and

$$\Sigma = E[(x - \mu)(x - \mu)^T],$$

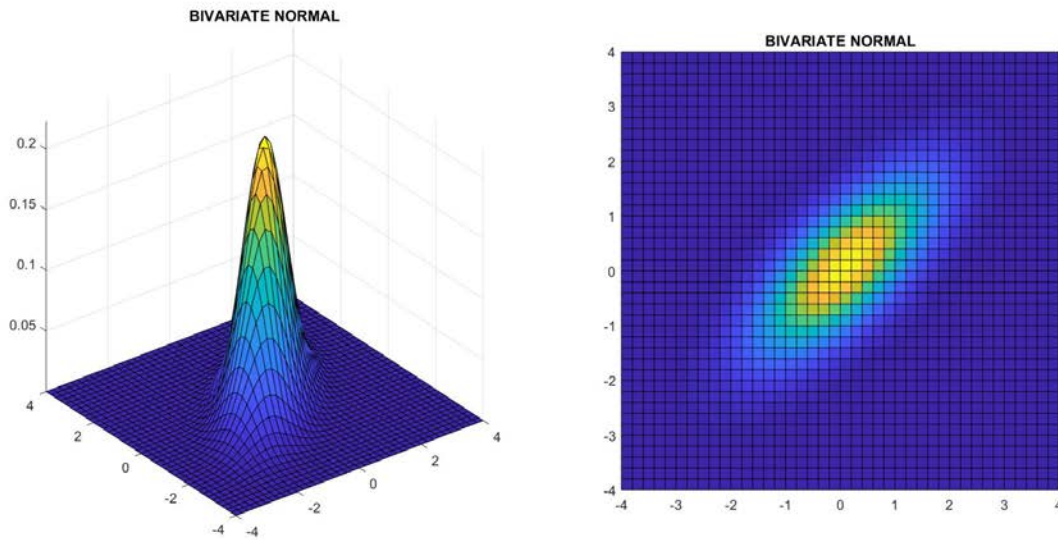
where the expected value of an array is given by the expected values of its components. The covariance matrix is symmetric ($\Sigma^T = \Sigma$) positive definite (all eigenvalues of Σ are greater than zero) for most applications of interest to statisticians and engineers.

We illustrate some properties of the multivariate normal by looking at the bivariate ($d = 2$) case. The probability density function for a bivariate normal is represented by a bell-shaped surface. The center of the surface is determined by the mean μ and the shape of the surface is determined by the covariance Σ .

- If the covariance matrix is diagonal (all of the off-diagonal elements are zero), and the diagonal elements are equal, then the shape is circular.
- If the diagonal elements are not equal, then we get an ellipse with the major axis vertical or horizontal. If the covariance matrix is not diagonal, then the shape is elliptical with the axes at an angle. Some of these possibilities are illustrated in the next example.



This figure shows a standard bivariate normal probability density function that is centered at the origin. The covariance matrix is given by the identity matrix. Notice that the shape of the surface looks circular. The plot on the right is for a viewpoint looking down on the surface.



This shows a bivariate normal density where the covariance matrix has non-zero off-diagonal elements. Note that the surface has an elliptical shape. The plot on the right is for a viewpoint looking down on the surface.

Generating Random Variables

Uniform Random Numbers

Most methods for generating random variables start with random numbers that are uniformly distributed on the interval $(0, 1)$. We will denote these random variables by the letter u . With the advent of computers, we now have the ability to generate uniform random variables very easily.

It should be noted that random numbers that are uniformly distributed over an interval a to b may be generated by a simple transformation, as follows

$$X = (b - a) \cdot u + a$$

where

$$u \sim U(0, 1) \text{ and } X \sim U(a, b).$$

Normal Random Variables

Given a standard normal random variable $Z \sim N(0, 1)$, we can obtain any normally distributed random variable X with mean μ and variance σ^2 by means of a transformation:

$$X = \mu + Z \cdot \sigma$$

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

Jointly Normal Random Variables

Suppose we generate standard normal random variables X and Y with correlation ρ

- Let X be standard normal
- Let U be standard normal (independent of X)
- Let $Y = \rho X + \sqrt{1 - \rho^2} U$
- $E(Y) = 0$, $Var(Y) = \rho^2 + 1 - \rho^2 = 1$,
- $Cov(X, Y) = E(XY) = \rho Var(X) = \rho$

So X and Y are standard normal with correlation ρ .

Inverse Transform Method

The inverse transform method can be used to generate random variables from a continuous distribution. It uses the fact that the cumulative distribution function F is uniform $(0, 1)$.

$$U = F(X).$$

If U is a uniform $(0, 1)$ random variable, then we can obtain the desired random variable X from the following relationship.

$$X = F^{-1}(U).$$

We see an example of how to use the inverse transform method when we discuss generating random variables from the exponential distribution (see the following Example). The general procedure for the inverse transformation method is outlined here.

PROCEDURE - INVERSE TRANSFORM METHOD (CONTINUOUS)

1. Derive the expression for the inverse distribution function $F^{-1}(U)$.
2. Generate a uniform random number U .
3. Obtain the desired X from $X = F^{-1}(U)$.

This same technique can be adapted to the discrete case. Say we would like to generate a discrete random variable X that has a probability mass function given by

$$P(X = x_i) = p_i; \quad x_1 < x_2 < \dots; \quad \sum_i p_i = 1.$$

We get the random variables by generating a random number U and then deliver the random number X according to the following

$$X = x_i \quad \text{if } F(x_{i-1}) < U \leq F(x_i).$$

Example: Exponential Distribution

The inverse transform method can be used to generate random variables from the exponential distribution and serves as an example of this procedure. The distribution function for an exponential random variable with parameter λ is given by

$$F(x) = 1 - e^{-\lambda x} \quad 0 < x < \infty.$$

Letting

$$u = F(x) = 1 - e^{-\lambda x},$$

we can solve for x , as follows

$$\begin{aligned} u &= 1 - e^{-\lambda x} \\ e^{-\lambda x} &= 1 - u \\ -\lambda x &= \log(1 - u) \\ x &= -\frac{1}{\lambda} \log(1 - u). \end{aligned}$$

By making note of the fact that $1 - u$ is also uniformly distributed over the interval $(0,1)$, we can generate exponential random variables with parameter λ using the transformation

$$X = -\frac{1}{\lambda} \log(U).$$

Acceptance-Rejection Method

In some cases, we might have a simple method for generating a random variable from one density, say $g(y)$, instead of the density we are seeking. We can use this density to generate from the desired continuous density $f(x)$. We first generate a random number Y from $g(y)$ and accept the value with a probability proportional to the ratio $f(Y)/g(Y)$.

If we define c as a constant that satisfies

$$\frac{f(y)}{g(y)} \leq c; \quad \text{for all } y,$$

then we can generate the desired variates using the procedure outlined below. The constant c is needed because we might have to adjust the height of $g(y)$ to ensure that it is above $f(y)$. We generate points from $cg(y)$, and those points that are inside the curve $f(y)$ are accepted as belonging to the desired density. Those that are outside are rejected. It is best to keep the number of rejected variates small for maximum efficiency.

1. Choose a density $g(y)$ that is easy to sample from.
2. Find a constant c such that Equation $\frac{f(y)}{g(y)} \leq c$ is satisfied.
3. Generate a random number Y from the density $g(y)$.
4. Generate a uniform random number U .
5. If

$$U \leq \frac{f(Y)}{cg(Y)},$$

then accept, else go to step 3.

Metropolis-Hastings Algorithms

The Metropolis-Hastings method is a generalization of the Metropolis technique of Metropolis, et al. [1953], which had been used for many years in the physics community. The paper by Hastings [1970] further generalized the technique in the context of statistics. The Metropolis sampler, the independence sampler and the random-walk are all special cases of the Metropolis-Hastings method. Thus, we cover the general method first, followed by the special cases.

These methods share several properties, but one of the more useful properties is that they can be used in applications where is known up to the constant of proportionality. Another property that makes them useful in a lot of applications is that the analyst does not have to know the conditional distributions, which is the case with the Gibbs sampler. While it can be shown that the Gibbs sampler is a special case of the Metropolis-Hastings algorithm. We include it in the next section because of this difference.

Metropolis-Hastings Sampler

The Metropolis-Hastings sampler obtains the state of the chain at by sampling a *candidate point* Y from a *proposal distribution* . Note that this depends only on the previous state and can have any form, subject to regularity conditions. An example for is the multivariate normal with mean and fixed covariance matrix. One thing to keep in mind when selecting is that the proposal distribution should be easy to sample from.

The required regularity conditions for are irreducibility and aperiodicity [Chib and Greenberg, 1995]. **Irreducibility** means that there is a positive probability that the Markov chain can reach any non-empty set from all starting points. **Aperiodicity** ensures that the chain will not oscillate between different sets of states. These conditions are usually satisfied if the proposal distribution has a positive density on the same support as the target distribution. They can also be satisfied when the target distribution has a restricted support. For example, one could use a uniform distribution around the current point in the chain.

The candidate point is accepted as the next state of the chain with probability given by

$$\alpha(X_t, Y) = \min\left\{1, \frac{\pi(Y)q(X_t|Y)}{\pi(X_t)q(Y|X_t)}\right\}.$$

If the point Y is not accepted, then the chain does not move and $X_{t+1} = X_t$. The steps of the algorithm are outlined below. It is important to note that the distribution of interest $\pi(x)$ appears as a ratio, so the constant of proportionality cancels out. This is one of the appealing characteristics of the Metropolis-Hastings sampler, making it appropriate for a wide variety of applications.

1. Initialize the chain to X_0 and set $t = 0$.
2. Generate a candidate point Y from $q(\cdot | X)$.
3. Generate U from a uniform $(0, 1)$ distribution.
4. If $U \leq \alpha(X_t, Y)$ then set $X_{t+1} = Y$, else set $X_{t+1} = X_t$.
5. Set $t = t + 1$ and repeat steps 2 through 5.

The Metropolis-Hastings procedure is implemented in Example, where we use it to generate random variables from a standard Cauchy distribution. As we will see, this implementation is one of the special cases of the Metropolis-Hastings sampler described later.

Example

We show how the Metropolis-Hastings sampler can be used to generate random variables from a standard Cauchy distribution given by

$$f(x) = \frac{1}{\pi(1 + x^2)} \quad -\infty < x < \infty.$$

From this, we see that

$$f(x) \propto \frac{1}{1 + x^2}$$

We will use the normal as our proposal distribution, with a mean given by the previous value in the chain and a standard deviation given by σ .

Non-parametric Kernel Density

A kernel distribution is a nonparametric representation of the probability density function (pdf) of a random variable. You can use a kernel distribution when a parametric distribution cannot properly describe the data, or when you want to avoid making assumptions about the distribution of the data. A kernel distribution is defined by a smoothing function and a bandwidth value, which control the smoothness of the resulting density curve.

Kernel Density Estimator

The kernel density estimator is the estimated pdf of a random variable. For any real values of x , the kernel density estimator's formula is given by

$$\hat{f}_h = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

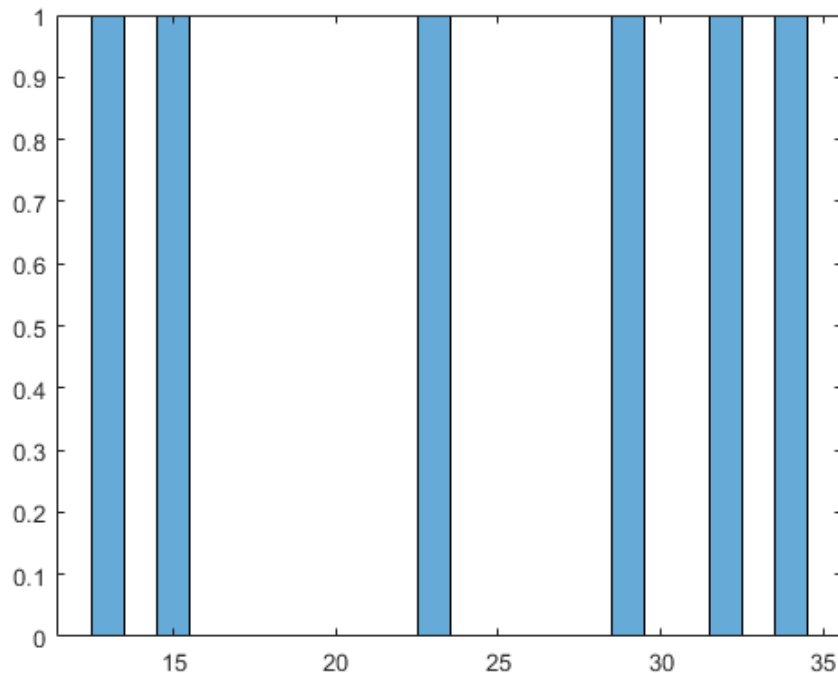
where x_1, x_2, \dots, x_n are random samples from an unknown distribution, n is the sample size, $K(\bullet)$ is the kernel smoothing function, and h is the bandwidth.

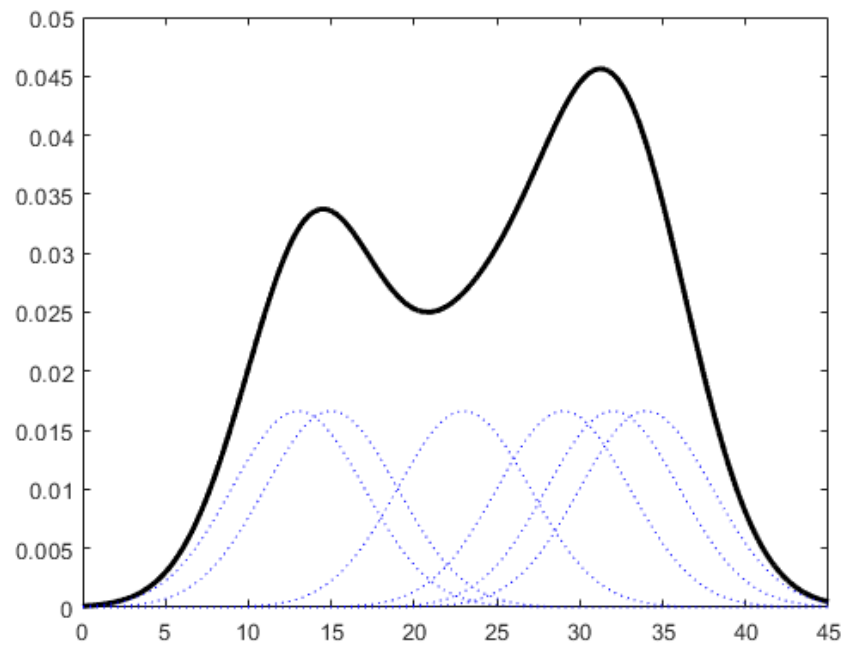
Kernel Smoothing Function

The kernel smoothing function defines the shape of the curve used to generate the pdf. Similar to a histogram, the kernel distribution builds a function to represent the probability distribution using the sample data. But unlike a histogram, which places the values into discrete bins, a kernel distribution sums the component smoothing functions for each data value to produce a smooth, continuous probability curve. You may plot a visual comparison of a histogram and a kernel distribution generated from the same sample data.

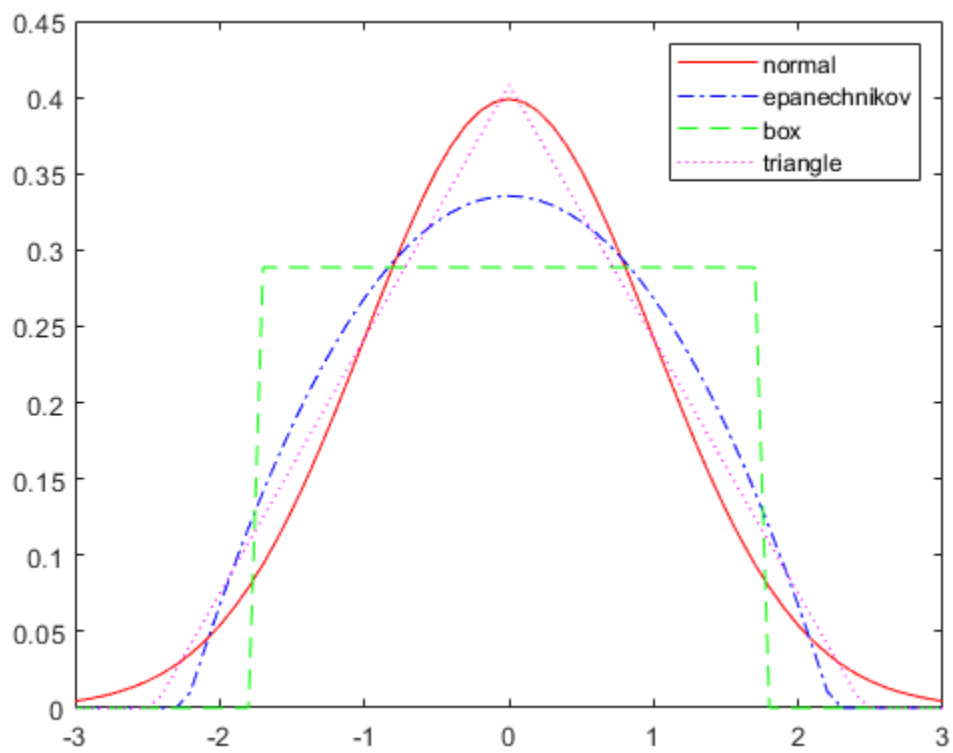
Because of this bin count approach, the histogram produces a discrete probability density function. This might be unsuitable for certain applications, such as generating random numbers from a fitted distribution.

Alternatively, the kernel distribution builds the pdf by creating an individual probability density curve for each data value, then summing the smooth curves. This approach creates one smooth, continuous probability density function for the data set.





Examples of Kernels for Density Estimation



- Triangle $K(t) = (1 - |t|) - 1 \leq t \leq 1$
- Epanechnikov

$$K(t) = \frac{3}{4}(1 - t^2) - 1 \leq t \leq 1$$

- Biweight

$$K(t) = \frac{15}{16}(1 - t^2)^2 - 1 \leq t \leq 1$$

- Triweight

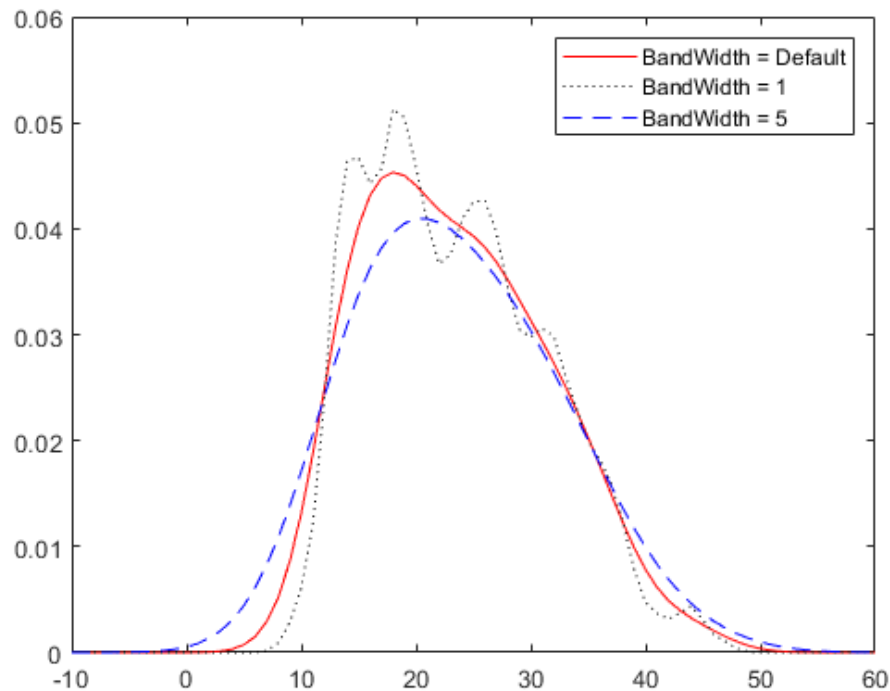
$$K(t) = \frac{35}{32}(1 - t^2)^3 - 1 \leq t \leq 1$$

- Normal

$$K(t) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{t^2}{2}\right\}$$

Bandwidth

The choice of bandwidth value controls the smoothness of the resulting probability density curve. This plot shows the density estimate for the same MPG data, using a normal kernel smoothing function with three different bandwidths.



Bootstrap

The bootstrap is a simulation method for forming confidence intervals and obtaining standard errors using only information from the sample. Like a Monte-Carlo simulation, but uses only the data.

Advantages

1. Applicable in a wide range of contexts
2. Easy.
3. The bootstrap in some cases may produce better approximations (knocks out an extra term in the Edgeworth expansion).

Hall justifies the bootstrap with a “Russian Dolls” analogy

- Doll zero: population that we do not get to see
- Doll one: sample we observe
- Doll two: bootstrap sample

Bootstrap: Sample Code in Matlab

```
X=[79 73 68 77 86 71 69]';
[T, N] = size(X);
x_mu = mean(X);
x_se = std(X)/sqrt(T);
B = 1000;
x_boot_mu = zeros(B,1);
for i = 1:B;
    x_boot=x(ceil(T*rand(T,1)+0.0001)); % draw x with replacement.
    x_boot_mu(i) = mean(x_boot);
    x_boot_se = std(x_boot)/sqrt(T);
    t_stat_boot(i)=(x_boot_mu(i)-x_mu)/x_boot_se;
end;

xbootmu = sort(x_boot_mu);
tstatboot=sort(t_stat_boot);

% confidence interval of x_boot_mu:
[x_boot_mu(25) x_boot_mu(975)]
% confidence interval of x:
[x_mu-(tstatboot(975)*x_se) x_mu-(tstatboot(25)*x_se)]
```

The last two lines give the OP and Percentile t CIs for mean

The bootstrap in a regression model

Suppose I have a linear regression model

$$y = \beta' x_i + \varepsilon_i$$

The most standard implementation of the bootstrap entails the following steps:

1. Estimate the parameter vector β and work out the residuals

$$e_i = y - \hat{\beta}' x_i$$

2. Resample from the residuals with replacement and from the regressors with replacement.
3. Build up a new dataset of the dependent variables as

$$y_i^{\text{BOOT}} = \hat{\beta}' x_i^{\text{BOOT}} + e_i^{\text{BOOT}}$$

4. Work out the quantity of interest in this new dataset
5. Repeat (2)-(4) many times.

Other percentile, percentile or percentile-t confidence intervals can then be worked out.

Disadvantage

- Destroy conditional Heteroskedasticity
- Destroy autocorrelation, as it assumes all observations are independent of each other.