

# Bivariate Probability Distributions

Abby Spurdle

April 11, 2021

*Convenience functions for plotting/evaluating bivariate (uniform, binomial, Poisson, categorical, normal and bimodal) distributions, trivariate (normal and Dirichlet) distributions, bivariate kernel density estimates and bivariate empirical cumulative distribution functions. Supports their probability mass functions (PMFs), probability density functions (PDFs) and cumulative distribution functions (CDFs), mostly where applicable.*

## Introduction

This package contains convenience functions for constructing, plotting and evaluating bivariate probability distributions, including their mass (PMF), density (PDF) and distribution (CDF) functions.

It supports the following parametric probability distributions:

- Discrete bivariate **Uniform** distributions (PMF and CDF).
- Bivariate **Binomial** distributions (PMF and CDF).
- Bivariate **Poisson** distributions (PMF and CDF).
- Bivariate **Categorical** distributions (PMF).
- Continuous bivariate **Uniform** distributions (PDF and CDF).
- Bivariate **Normal** distributions (PDF and CDF).
- Trivariate **Normal** distributions (PDF).
- Bivariate **Bimodal** distributions (PDF and CDF).
- Trivariate **Dirichlet** distributions (PDF).

And it supports the following nonparametric probability distributions:

- Bivariate **Kernel** density estimates (PDF).
- Bivariate **Empirical** cumulative distribution functions (CDF).

There's more than one way of formulating bivariate binomial, Poisson and bimodal distributions. I've used the simplest approaches that I could for the binomial and bimodal distributions, with the Poisson distribution adapted from Karlis and Ntzoufras (2003).

Some of these distributions (color-coded in gold, or brown) are equivalent to the product of their marginal distributions. Others (color-coded in blue) may be equivalent to the product of their marginal distributions in some cases, but are not, in generality.

(i.e. A bivariate normal distribution can be represented as the product of two univariate normal distributions if it has no correlation).

Note that:

- This package uses a system of self-referencing S4-based function objects. This allows us to plot and evaluate functions, without specifying their parameters, each time.
- The functions for evaluating discrete probability distributions, coerce their arguments to integers. If you try to evaluate discrete probability distributions with non-integer arguments, you may get unexpected results.
- The probhat package provides more tools for categorical distributions and kernel smoothing.

Refer to the barsurf package for information on how to customize plots, if required. (Plotting functions in this package call plotting functions from the barsurf package).

Also note that this vignette contains a small amount non-visible R code to change plot margins.

## Preliminary Code (And Required Packages)

I will load (and attach) the barsurf and bivariate packages:

```
> library (barsurf)
> library (bivariate)
```

Note that the bivariate package imports the barsurf, mvtnorm and KernSmooth packages, and the barsurf package imports the kubik and colorspace packages.

Also, the misc3d package needs to be installed and loaded, in order to plot the trivariate normal distribution.

And I will use a dataset from the MASS package, later.

I will set global options:

```
> set.bs.options (nhl=0, rendering.style="pdf", theme="gold", ref.arrows=FALSE)
```

The “pdf” rendering style uses finer lines.

Also, I’m setting the theme (for default colors) to gold, but I will use blue for many plots.

## Discrete **Bivariate** Uniform Distributions

We can describe a discrete bivariate uniform distribution as the product of two discrete univariate uniform distributions, so:

$$\begin{aligned} \mathbb{P}(X = x, Y = y) &= f_{X,Y}(x, y; a_X, b_X, a_Y, b_Y) \\ &= f_X(x; a_X, b_X) f_Y(y; a_Y, b_Y) \quad (\text{here, } x \text{ and } y \text{ are ignored}) \\ &= \left[ \frac{1}{b_X - a_X + 1} \right] \left[ \frac{1}{b_Y - a_Y + 1} \right] \end{aligned}$$

Where  $a_X$  and  $b_X$  are integers giving the minimum and maximum possible values of  $X$ , and  $a_Y$  and  $b_Y$  are the same, but for  $Y$ .

And assuming that  $x$  is in the interval  $[a_X, b_X]$  and  $y$  is in the interval  $[a_Y, b_Y]$ .

We can construct its probability mass function using the **dubvpmf** function, and its cumulative distribution function using the **dubvcdf** function.

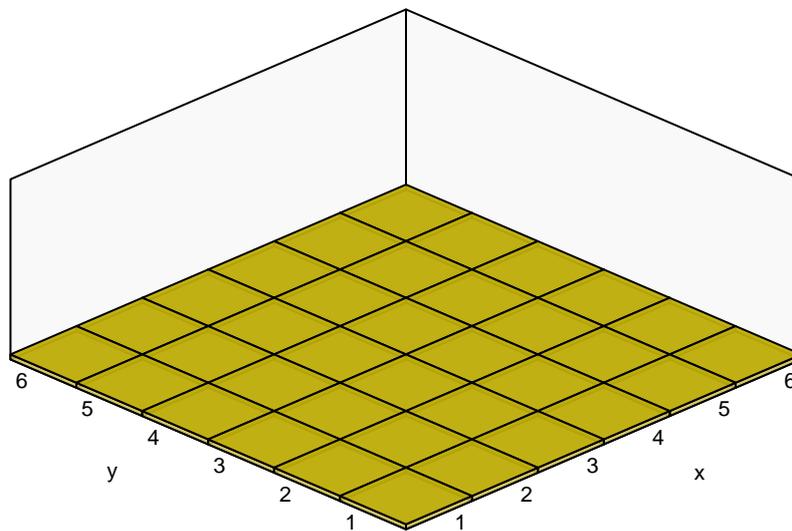
Both constructors take four arguments.

Here's an example, where both  $X$  and  $Y$ , can take values between one and six:

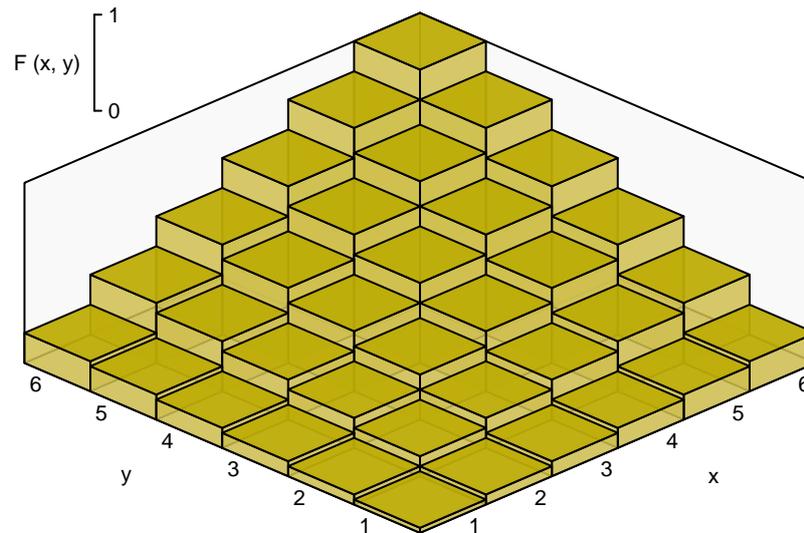
```
> f <- dubvpmf (  
  1, 6, #first variable  
  1, 6) #second variable  
> F <- dubvcdf (  
  1, 6,  
  1, 6)
```

And we can plot the functions:

```
> plot (f)
```



```
> plot (F)
```



Reiterating, this package uses a system of self-referencing function objects, which allows us to plot and evaluate functions, without specifying their parameters, each time.

In both cases, we can evaluate the functions for  $x$  and  $y$ :

```
> f (2, 4)
[1] 0.02777778
> F (2, 4)
[1] 0.2222222
```

The same applies to all the other probability distributions in this package, except for kernel density estimates.

Note that we can compute these values from univariate distributions: (Again, assuming that  $x$  and  $y$  are within the supported region).

```
> d.unif.pmf.eval <- function (x, a, b)    #x ignored
  1 / (b - a + 1)
> d.unif.cdf.eval <- function (x, a, b)    #x used
  (x - a + 1) / (b - a + 1)

> d.unif.pmf.eval (2, 1, 6) * d.unif.pmf.eval (4, 1, 6)
[1] 0.02777778
> d.unif.cdf.eval (2, 1, 6) * d.unif.cdf.eval (4, 1, 6)
[1] 0.2222222
```

This only applies to the probability distributions color-coded in gold.

## Bivariate **Binomial** Distributions

One way to define a bivariate binomial distribution is to say that we have  $n$  trials. In each trial there are two independent events, each with a particular probability of success. Like flipping two coins,  $n$  times.

Like the bivariate uniform distribution, we can describe a bivariate binomial distribution

as the product of two univariate binomial distributions, with the same  $n$  parameter, so:

$$\begin{aligned}\mathbb{P}(X = x, Y = y) &= f_{X,Y}(x, y; p_X, p_Y, n) \\ &= f_X(x; p_X, n) f_Y(y; p_Y, n) \\ &= \left[ \binom{n}{x} p_X^x (1 - p_X)^{n-x} \right] \left[ \binom{n}{y} p_Y^y (1 - p_Y)^{n-y} \right]\end{aligned}$$

Where  $p_X$  is the probability of the first success and  $p_Y$  is the probability of the second success.

We can construct its probability mass function using the `bnbvpmf` function, and its cumulative distribution function using the `bnbvCDF` function.

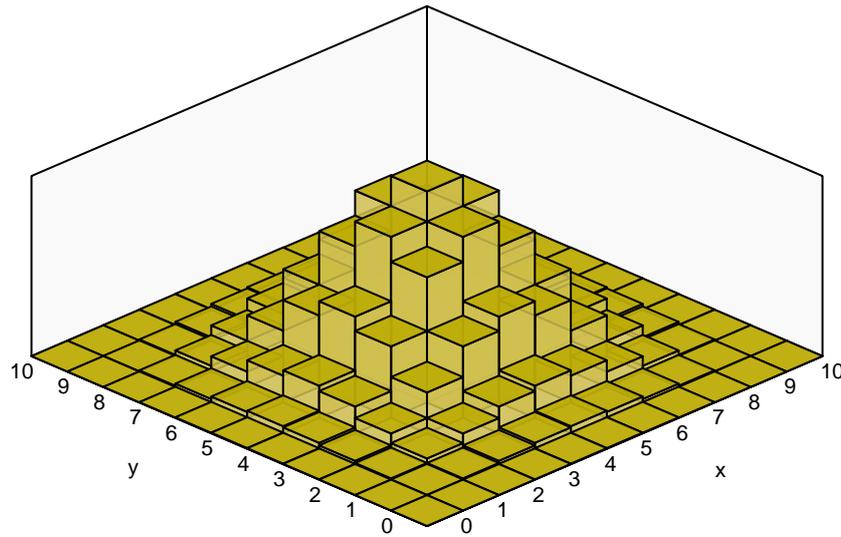
It takes three arguments, the probability of the first success, the probability of the second success and the number of trials.

Here's an example where the probability of both the first and second success is 0.5, and there's ten trials:

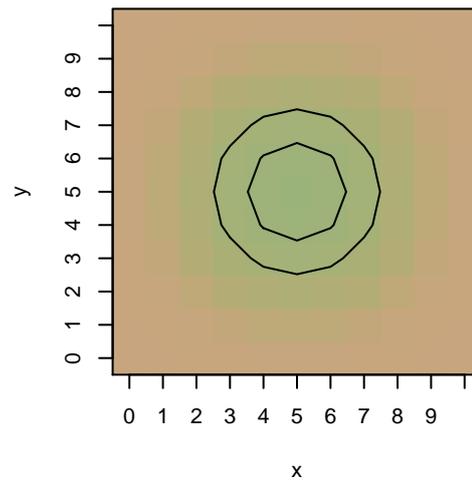
```
> f <- bnbvpmf (0.5, 0.5, 10)
> F <- bnbvCDF (0.5, 0.5, 10)
```

And again we can plot the functions:

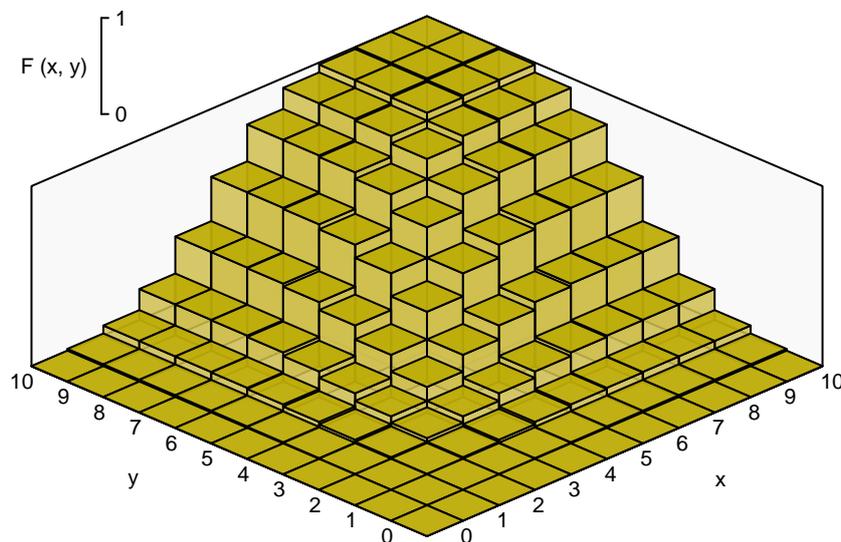
```
> plot (f)
```



```
> plot (f, FALSE)
```



```
> plot (F)
```



Note that I've put the probabilities ("p") first, however, it's customary for the number of trials ("n") to go first.

If n is omitted, it defaults to one, giving a bivariate Bernoulli distribution.

## Bivariate Poisson Distributions

Based on Karlis and Ntzoufras (2003), we can define the bivariate Poisson probability mass function as:

$$\begin{aligned} \mathbb{P}(X = x, Y = y) &= f_{X,Y}(x, y; \lambda_1, \lambda_2, \lambda_3) \\ &= e^{-(\lambda_1 + \lambda_2 + \lambda_3)} \frac{\lambda_1^x \lambda_2^y}{x! y!} \sum_k \binom{x}{k} \binom{y}{k} k! \left( \frac{\lambda_3}{\lambda_1 \lambda_2} \right)^k \end{aligned}$$

Where:

$\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are positive real numbers.  
 $k$  is an integer in the sequence 0 to  $\min(x, y)$ .

And where:

$$\mathbb{E}(X) = \text{var}(X) = \lambda_1 + \lambda_3$$

$$\mathbb{E}(Y) = \text{var}(Y) = \lambda_2 + \lambda_3$$

$$\text{cov}(X, Y) = \lambda_3$$

Unlike the two previous probability distributions, this probability distribution is not the product of two marginal distributions.

We can construct its probability mass function using the `pbvpmf` or `pbvpmf.2` functions, and its cumulative distribution function using the `pbvcdf` or `pbvcdf.2` functions.

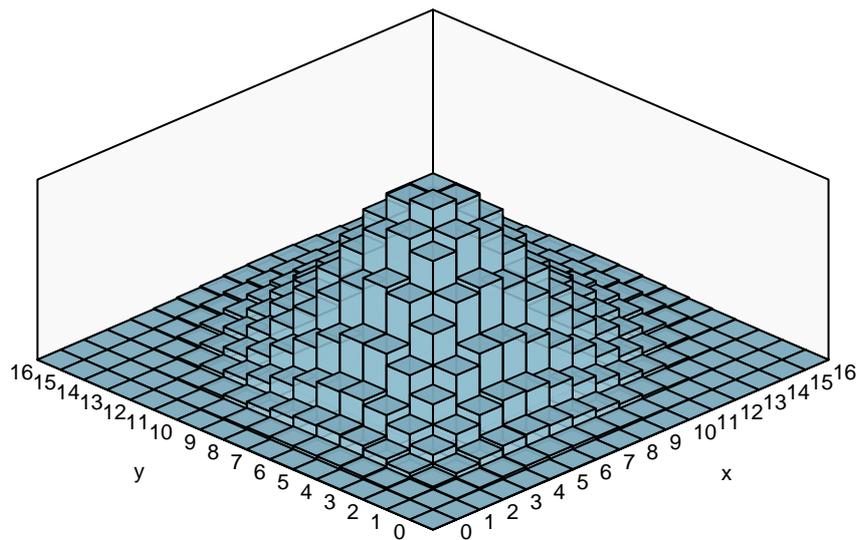
All functions take three arguments, the first versions (with no suffix) take the three  $\lambda$  parameters, and the second versions (with a suffix) take the expected value of  $X$ , the expected value of  $Y$  and the covariance between  $X$  and  $Y$ .

Here's an example where the expected value of both  $X$  and  $Y$  is eight, and the covariance is two:

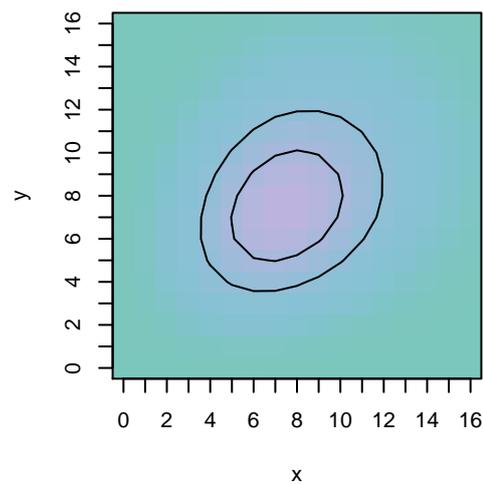
```
> f <- pbvpmf.2 (8, 8, 2)
```

And plots, in 2d and 3d:

```
> plot (f, theme="blue")
```



```
> plot (f, FALSE, theme="blue")
```



## Bivariate Categorical Distributions

We can construct a probability mass function for a categorical distribution, using the `gbvpmf` function.

It takes a single argument, a matrix of probabilities (or frequencies), preferably with row and column names.

Note that increasing row indices correspond to increasing x values and increasing column indices correspond to increasing y values.

Here's an example using a matrix of random counts:

```
> h <- matrix (sample (1:24), 4, 6)
> rownames (h) <- LETTERS [1:4]
> colnames (h) <- letters [1:6]

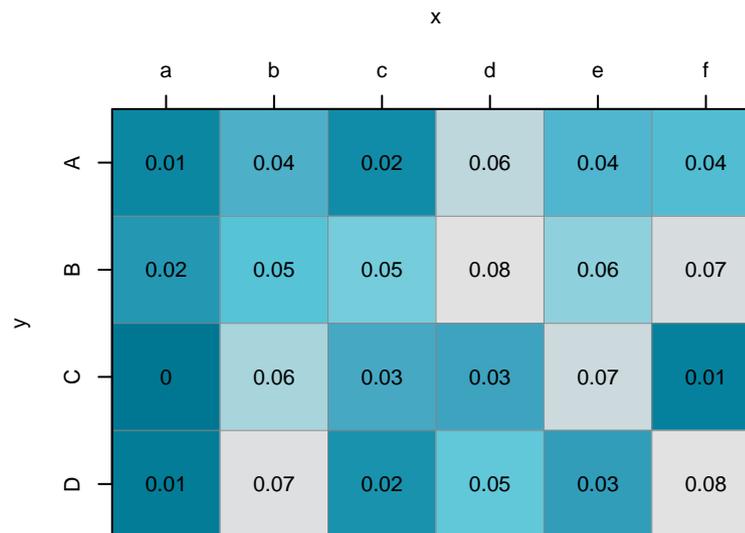
> h

   a  b  c  d  e  f
A  4 11  5 19 12 13
B  7 14 16 23 17 21
C  1 18 10  9 20  3
D  2 22  6 15  8 24

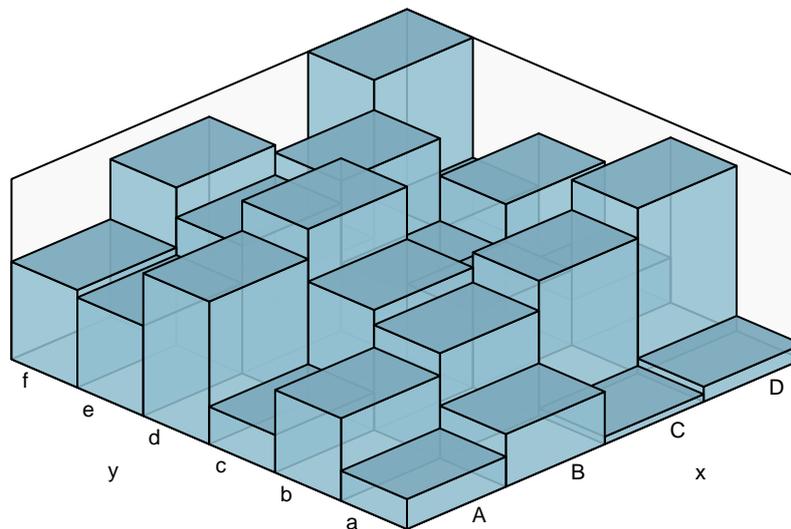
> f <- gbvpmf (h)
```

And plots in 2d and 3d:

```
> plot (f, FALSE, theme="blue")
```



```
> plot (f, theme="blue")
```



Note that the 2d plot has a different orientation from other 2d plots in this package.

Evaluation can use either integers or strings, and returns probabilities:

```
> h [2, 4] / sum (h)
```

```
[1] 0.07666667
```

```
> f (2, 4)
```

```
[1] 0.07666667
```

```
> f ("B", "d")
```

```
[1] 0.07666667
```

Note that the probhat package supports categorical distributions, and provides more tools.

## Continuous Bivariate Uniform Distributions

Continuous bivariate uniform distributions are similar to discrete bivariate uniform distributions. However, we have a probability density function rather than a probability mass function.

We can construct its probability density function using the  `cubvpdf`  function, and its cumulative distribution function using the  `cubvcdf`  function.

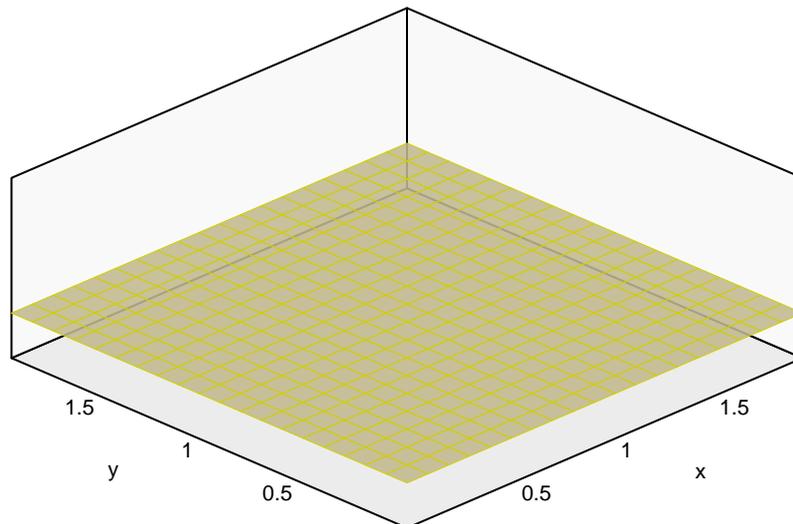
Both take four arguments.

Here's an example, where both  $X$  and  $Y$ , can take values between zero and two:

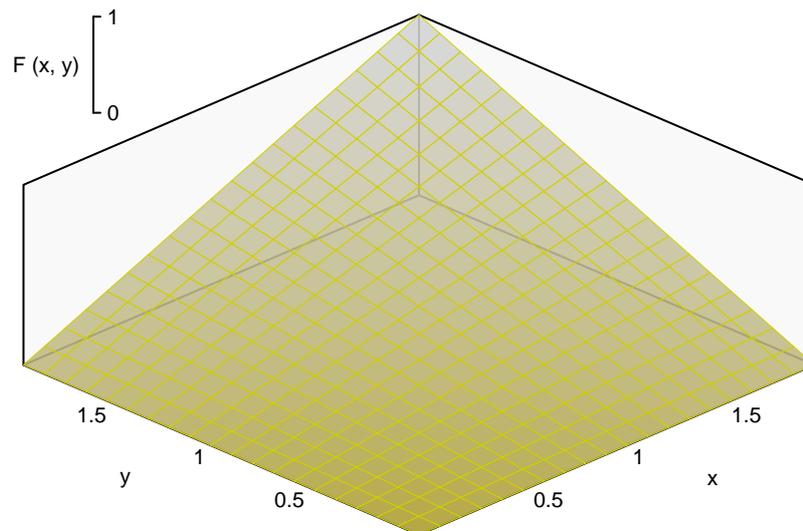
```
> f <- cubvpdf (
  0, 2, #first variable
  0, 2) #second variable
> F <- cubvcdf (
  0, 2,
  0, 2)
```

And plots of the functions:

```
> plot (f)
```



```
> plot (F)
```



## Bivariate Normal Distributions

This package uses the `mvtnorm` package to evaluate bivariate normal distributions. Please refer to that package for technical details.

We can construct a probability density function for the bivariate normal distribution using the `nbvpdf` or `nbvpdf.2` functions, and its cumulative distribution function using the `nbvcdf` or `nbvcdf.2` functions.

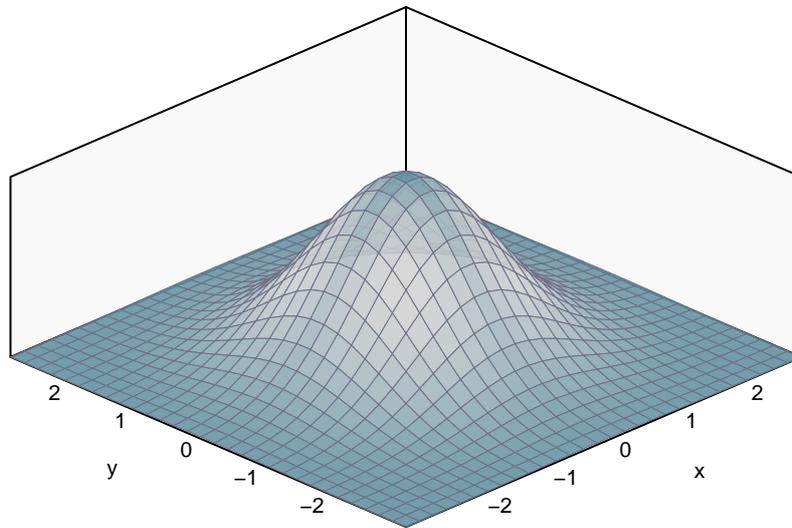
All functions take five parameters. The first functions (with no suffix) take the means of  $X$  and  $Y$ , the standard deviations of  $X$  and  $Y$ , and their correlation. The second functions (with the suffix) take the means of  $X$  and  $Y$ , the variances of  $X$  and  $Y$ , and their covariance.

Here's an example with zero means, standard deviations of one, and no correlation: (Essentially, a bivariate generalization of the "standard normal distribution").

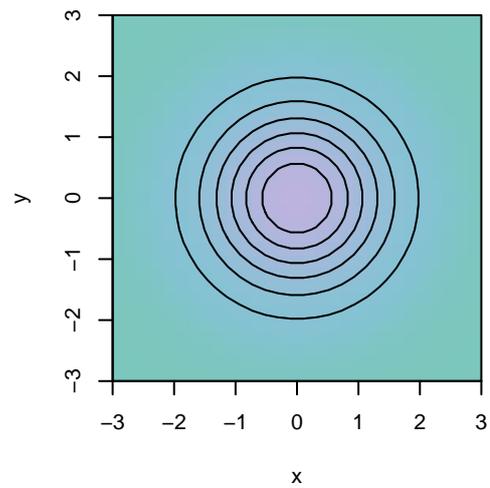
```
> f <- nbvpdf (0, 0, 1, 1, 0)
> F <- nbvcdf (0, 0, 1, 1, 0)
```

And plots in 2d and 3d:

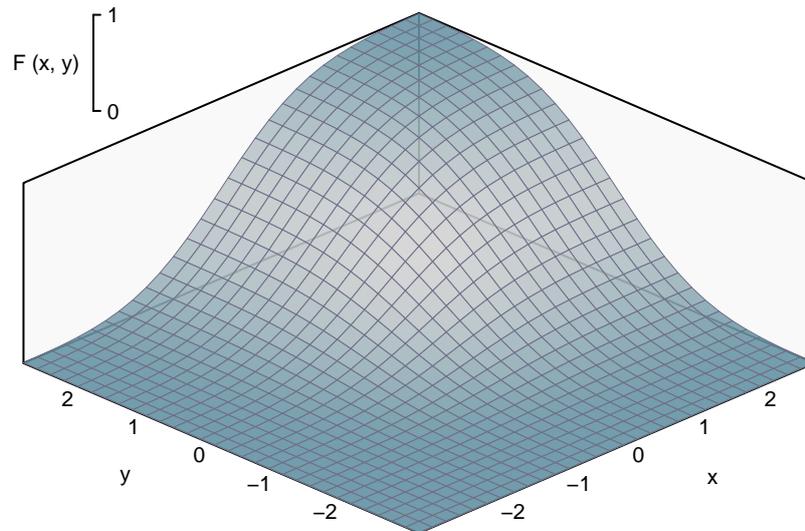
```
> plot (f, theme="blue")
```



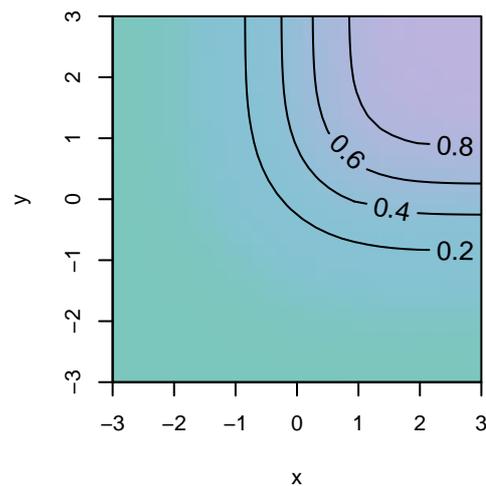
```
> plot (f, FALSE, theme="blue")
```



```
> plot (F, theme="blue")
```



```
> plot (F, FALSE, theme="blue")
```



Note that Appendix A shows the trivariate normal, and Appendix B compares bivariate normal distributions with different correlation parameters.

## Bivariate Bimodal Distributions

It's possible to construct a bivariate bimodal probability density function by taking two bivariate normal probability density functions, then adding their densities together, and then dividing by two.

We can construct such a probability density function using the `bmbvpdf` or `bmbvpdf.2` functions, and its cumulative distribution function using `bmbvcdf` or `bmbvcdf.2` functions. All functions take eight arguments, and follow the same principles as the normal distributions, discussed in the previous section. The first four arguments are the means and standard deviations (or variances) of the first component distribution. The last four ar-

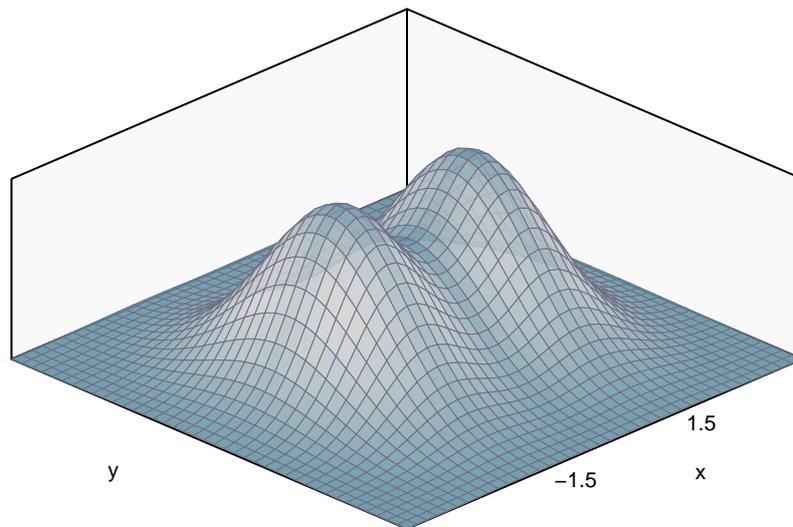
guments are the means and standard deviations (or variances) of the second component distribution.

Here's a an example, where the component means of  $X$  are 3.5 and 6.5, the component means of  $Y$  are zero, and all standard deviations are one:

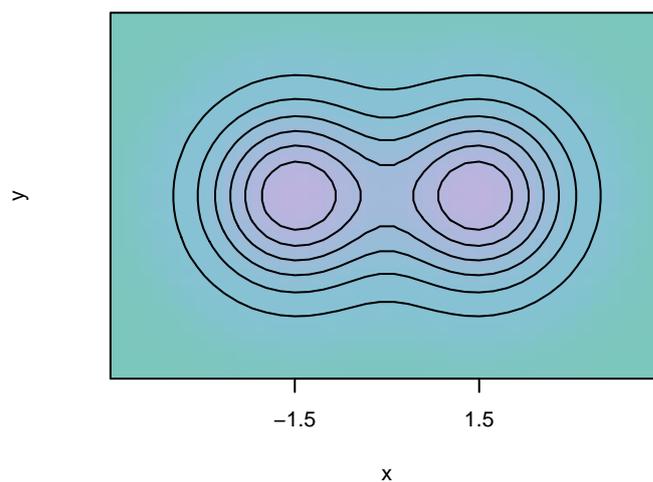
```
> f <- bmbvpdf (-1.5, 1.5)
```

And plots in 2d and 3d:

```
> plot (f, theme="blue",  
       axes = c (TRUE, FALSE), xat = c (-1.5, 1.5) )
```



```
> plot (f, FALSE, theme="blue",  
       axes = c (TRUE, FALSE), xat = c (-1.5, 1.5) )
```



Note that this method (of adding normal densities) is similar kernel smoothing, discussed later.

## Trivariate Dirichlet Distributions

Dirichlet distributions with three variables are similar to other probability distributions with two variables.

(Because it's possible to compute the third variable from the first two variables).

We can construct its probability density function using the `dtvpdf` function.

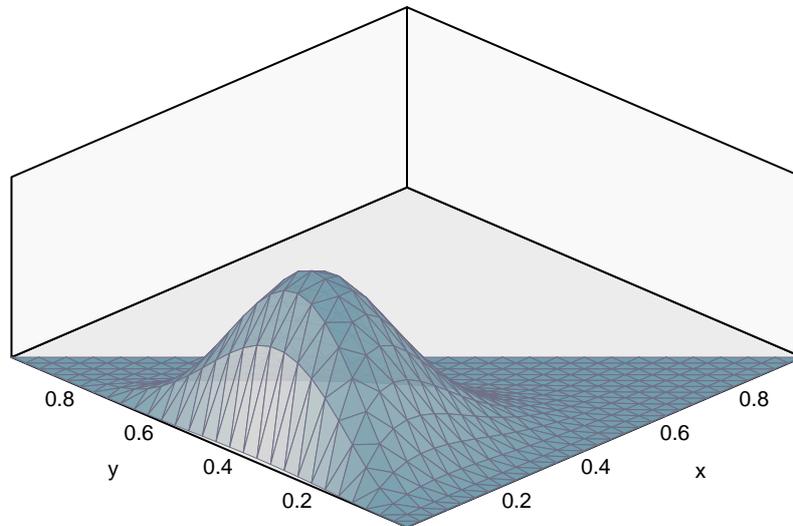
It takes three  $\alpha$  parameters.

Here's an example with  $\alpha$  parameters of two, four and six:

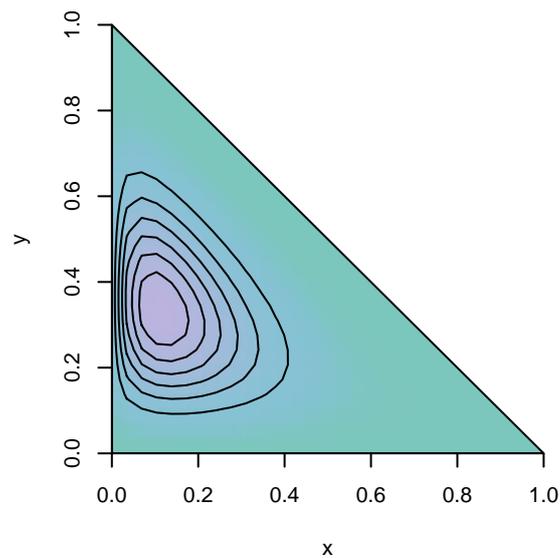
```
> f <- dtvpdf (2, 4, 6)
```

And plots in 2d and 3d:

```
> plot (f, theme="blue")
```



```
> plot (f, FALSE, theme="blue")
```



## Bivariate Kernel Density Estimates

This package uses the KernSmooth package to produce bivariate kernel density estimates. Please refer to that package for technical details.

We can construct a probability distribution representing bivariate kernel density estimates using the **kbvpdf** function.

It takes four arguments, two equal length vectors of data, and two bandwidth parameters.

I've adapted this example from KernSmooth:

(Which has a bandwidth parameter of 0.7 for duration, and 7 for waiting).

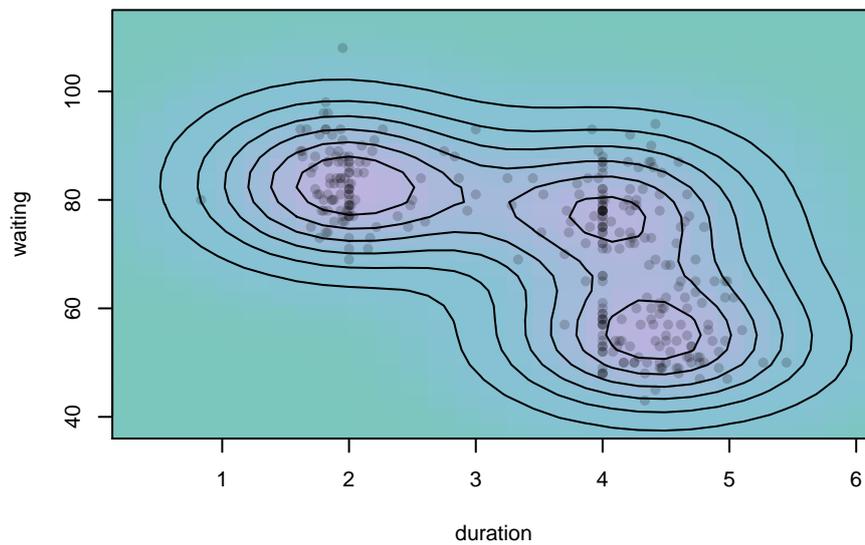
```
> data ("geyser", package="MASS")
> names (geyser)

[1] "waiting" "duration"

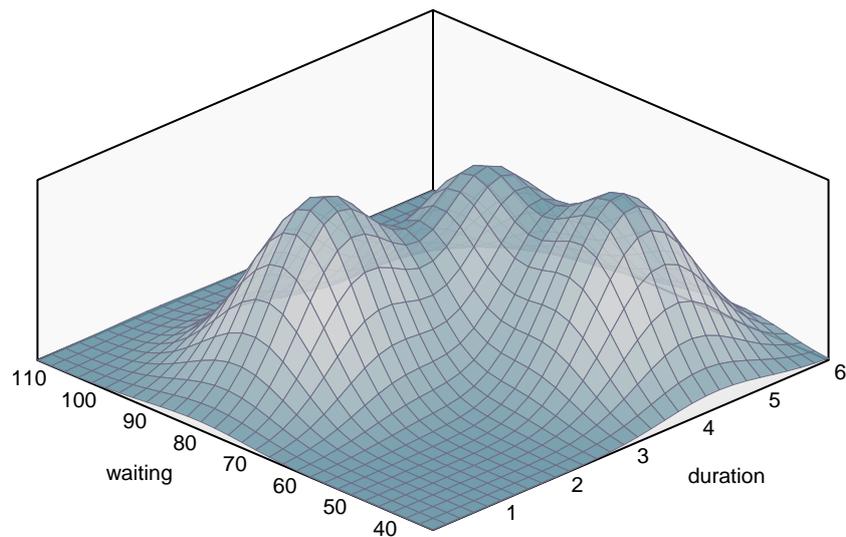
> fh <- kbvpdf (, ,0.7, 7, data = geyser [,2:1])
```

Again, once we have constructed our object we can plot it:

```
> plot (fh, FALSE, theme="blue")
```



```
> plot (fh, theme="blue")
```



Unlike other probability distributions in this package, you can't evaluate the function, `fh`. (However, the `bvmat` function can be used to compute density matrices).

Note that the `probhat` package provides more tools for kernel smoothing.

## Bivariate **Empirical** Cumulative Distribution Functions

Like bivariate kernel density estimates, bivariate ECDFs are computed from vectors of data.

The resulting probability distributions are step functions (so, have discrete properties), however in general, they represent the distributions of continuous random variables.

We can construct them using the `ebvcdf` function.

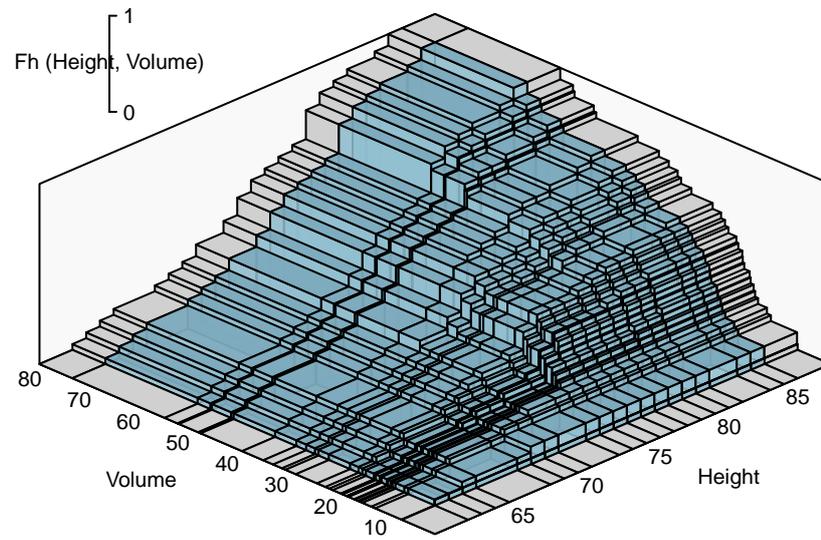
It takes two arguments, two equal length vectors of observations.

Here's an example, using the `trees` data:

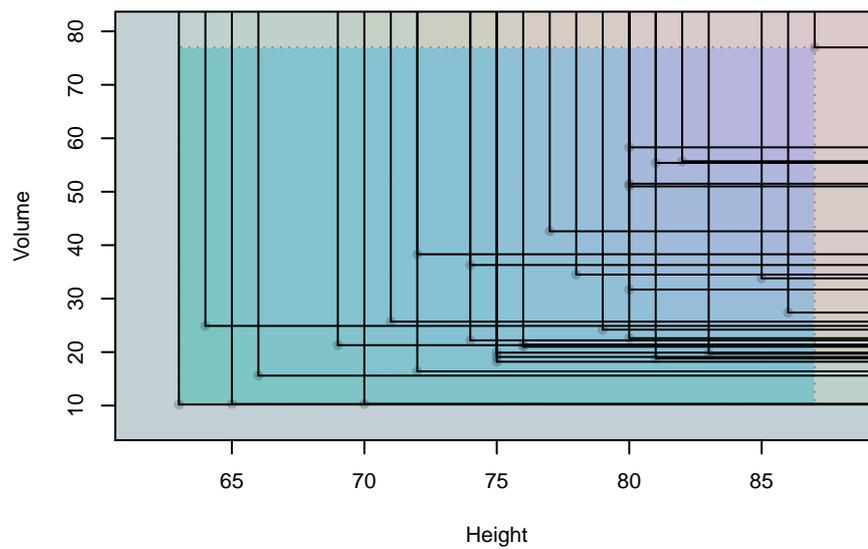
```
> names (trees)
[1] "Girth" "Height" "Volume"
> Fh <- ebvcdf (data = trees [,2:3])
```

And plots in 2d and 3d:

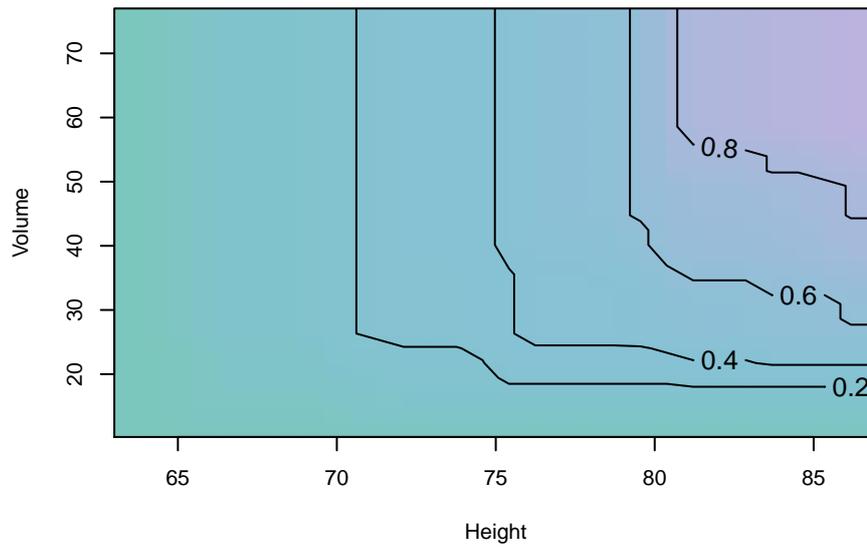
```
> plot (Fh, theme="blue")
```



```
> plot (Fh, FALSE, theme="blue")
```



```
> plot (Fh, FALSE, reg=TRUE, theme="blue")
```



Note that by default, if the number of observations is less than or equal to forty, the function is plotted as a (discrete) step function, with an extrapolated region outside the range of observed values. However, if the number of observations is larger than forty, the function is plotted as a (continuous) surface, evaluated over a regularly-spaced grid.

## References

### Imported R Packages

Spurdle, A. (2020). `barsurf`: Contour Plots, 3D Plots, Vector Fields and Heatmaps.

Genz, A., Bretz, F., Miwa, T., Mi, X. & Hothorn, T. (2020). `mvtnorm`: Multivariate Normal and t Distributions.

Wand, M. (2020). `KernSmooth`: Functions for Kernel Smoothing Supporting Wand & Jones (1995).

### Other R Packages

Feng, D., & Tierney, L. (2020) `misc3d`: Miscellaneous 3D Plots.

Spurdle, A. (2021). `probhat`: Multivariate Generalized Kernel Smoothing and Related Statistical Methods.

Spurdle, A. (2020). `kubik`: Cubic Hermite Splines and Related Optimization Methods.

Ihaka, R., Murrell, P., Hornik, K., Fisher, J. Stauffer, R., Wilke, C., McWhite, C., & Zeileis, A. (2020). `colorspace`: A Toolbox for Manipulating and Assessing Colors and Palettes.

Ripley, B. (2021). `MASS`: Support Functions and Datasets for Venables and Ripley's `MASS`.

### Journal Articles

Karlis, D. & Ntzoufras, I. (2003). Analysis of sports data by using bivariate Poisson models.

## Appendix A: Trivariate Normal Distributions

**Re-iterating, the misc3d package needs to be installed and loaded, in order to plot the trivariate normal distribution.**

Like the bivariate normal, this package uses the mvtnorm package to evaluate trivariate normal distributions.

We can construct a probability density function for the trivariate normal distribution using the `ntvpdf` or `ntvpdf.2` functions.

Both functions take nine parameters. The first function (with no suffix) take the means of  $X$ ,  $Y$  and  $Z$ , the standard deviations of  $X$ ,  $Y$  and  $Z$ , and their correlations  $\text{cor}(X, Y)$ ,  $\text{cor}(X, Z)$  and  $\text{cor}(Y, Z)$ .

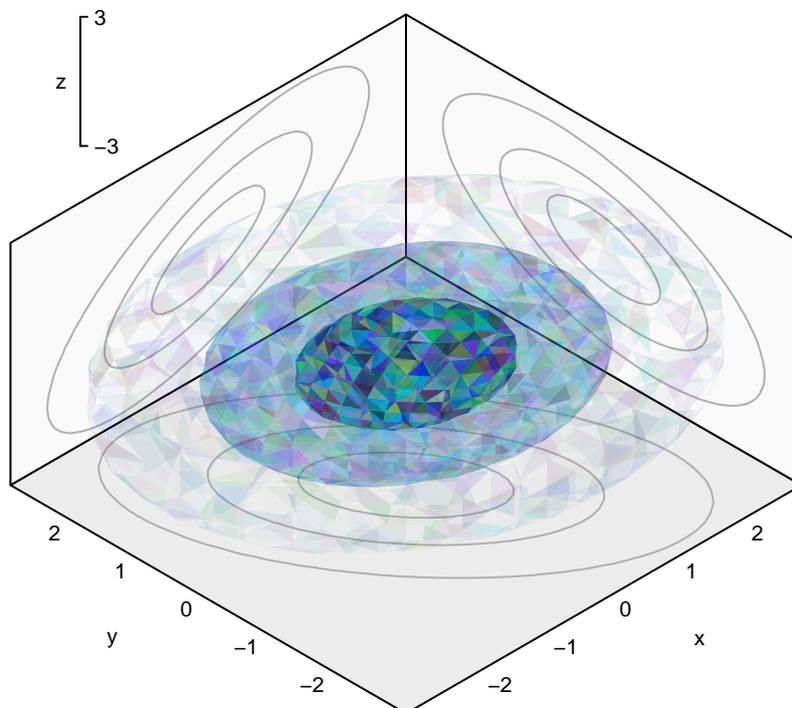
Here's an example with different correlation coefficients:

```
> f <- ntvpdf (
  0, 0, 0,      #mean: X, Y, Z
  1, 1, 1,      #sd:  X, Y, Z
  -0.5, 0.5, 0) #cor:  X~Y, X~Z, Y~Z

> f@covariance.matrix

      X   Y   Z
X  1.0 -0.5 0.5
Y -0.5  1.0 0.0
Z  0.5  0.0 1.0

> plot (f, theme="blue")
```



Note that the panel contours represent trivariate density values, where one variable (per

panel) is held constant.

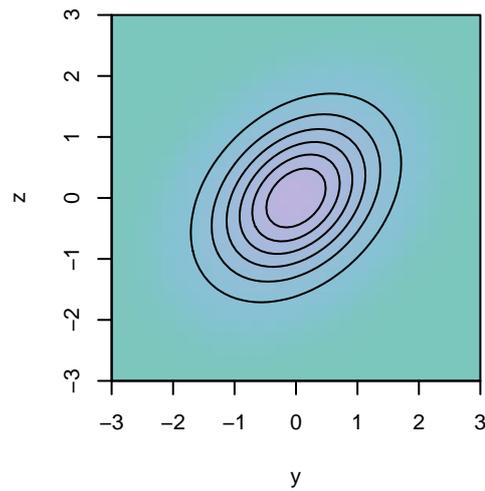
They do not represent the bivariate-marginal distributions.

Also, note the orientation of the plot.

(Increasing  $y$  values, run diagonally from bottom-right to top-left).

The following plot, gives the trivariate density values, where  $x$  is held constant at zero:

```
> plotf_cfield (function (y, z) f (0, y, z), c (-3, 3),
  theme="blue",
  xlab="y", ylab="z")
```



This is equivalent to the right-rear panel from the isosurface plot.

However, the horizontal axis has the reverse orientation to the  $y$ -axis from the isosurface plot.

In the trivariate distribution, there's zero correlation between  $Y$  and  $Z$ .

However, in the conditional distribution, there's positive correlation between  $Y$  and  $Z$ .

## Appendix B: Comparing Normal Distributions

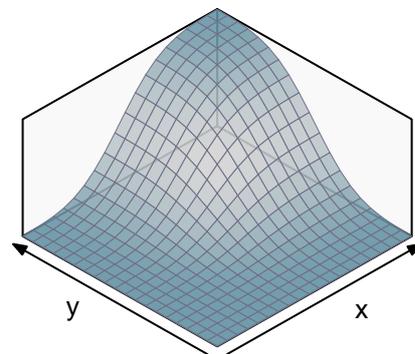
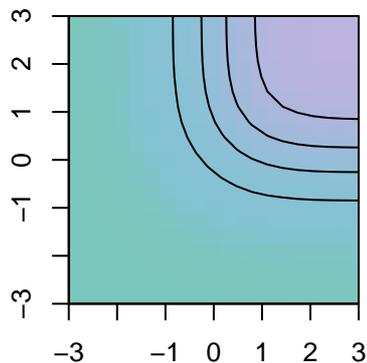
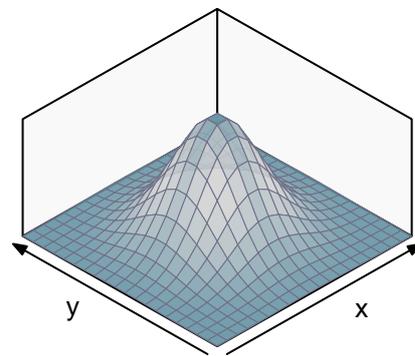
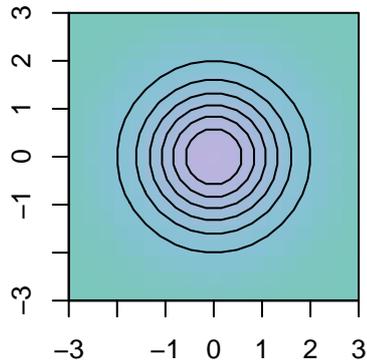
We can compare different normal distributions using different correlation or covariance parameters.

First, let's consider the bivariate distribution from the earlier section with zero correlation:

```
> f1 <- nbvpdf (0, 0, 1, 1, 0)
> f1@covariance.matrix

  X Y
X 1 0
Y 0 1

> plot (f1, all=TRUE, n=20, theme="blue")
```



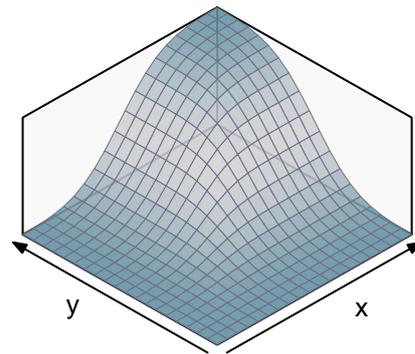
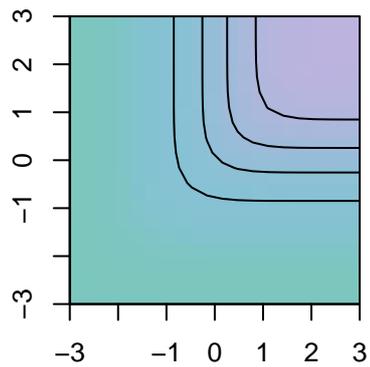
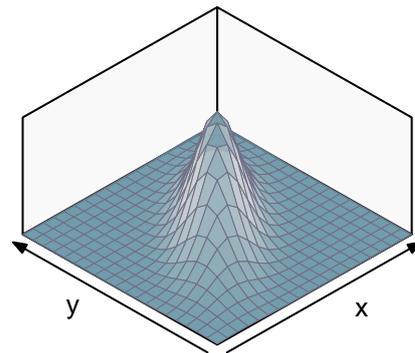
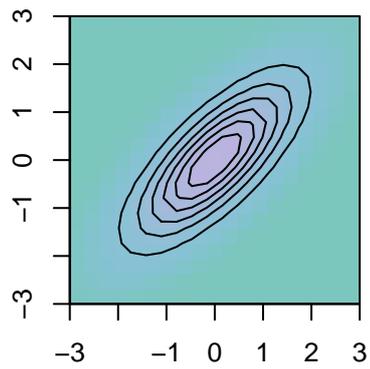
Second, let's consider bivariate distributions with positive correlation:

(Note that if both the standard deviations are one, then the correlation will equal the covariance).

```
> f2 <- nbvpdf (0, 0, 1, 1, 0.75)
> f2@covariance.matrix

  X   Y
X 1.00 0.75
Y 0.75 1.00
```

```
> plot (f2, all=TRUE, n=20, theme="blue")
```

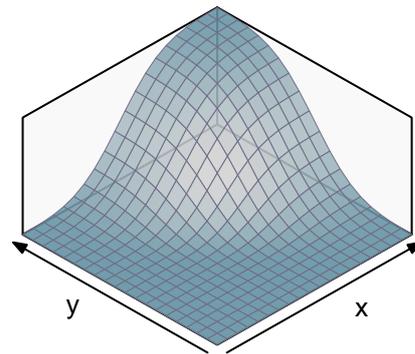
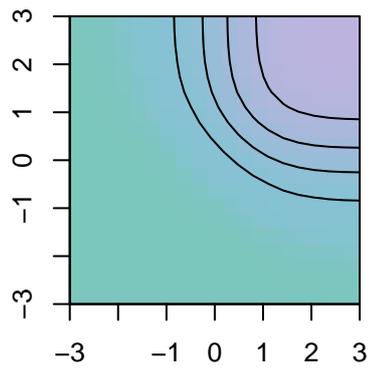
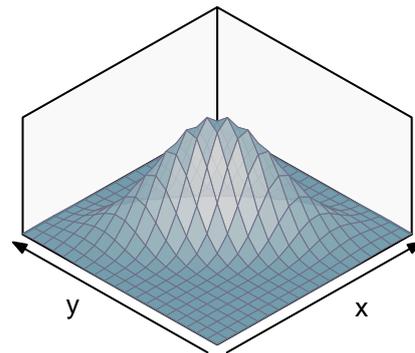
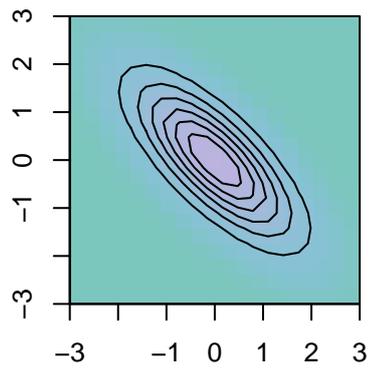


Third, let's consider bivariate distributions with negative correlation:

```
> f3 <- nbvpdf (0, 0, 1, 1, -0.75)  
> f3@covariance.matrix
```

	X	Y
X	1.00	-0.75
Y	-0.75	1.00

```
> plot (f3, all=TRUE, n=20, theme="blue")
```



Note that currently, the `all=TRUE` option requires the PMF or PDF rather than the CDF.