

Plotting an Invariant Mass Histogram in R

N.B.: Original text and code by [Edith Villegas Garcia](#). Edited and converted to a Jupyter notebook by [Achintya Rao](#). This notebook can be run in Binder by following the instructions from <https://github.com/binder-examples/dockerfile-r>.

If you're comfortable with R, [jump straight into the analysis of CMS Open Data](#). If not, let's start with a...

Quick intro to R

R is a language used for statistics and data analysis. In this tutorial, we introduce data analysis in R using the CMS Open Data from 2011. At the end of this tutorial, you will know how to create histograms of invariant mass and identify the particles that appear in the histograms.

Data types in R

There are different data types in R. Basic data types include the following:

- Logical
- Numeric
- Integer
- Complex
- Character

Variables are assigned to a type without having to specify it. Numeric variables are simply numbers. Integer variables are specified by adding an “L” at the end of the number:

```
[1] a <- 5L
```

Complex variables are complex numbers, they can be declared using the following format:

```
[2] b <- 5 + 3i  
d <- 8 + 0i
```

Logical variables can have values of either TRUE or FALSE. You can also assign a condition to a logical variable:

```
[3] c <- 3 > 5
```

Character variables are either letters or phrases, you can also add numbers and symbols between quotation marks.

```
[4] cr <- "3!"
```

Vectors

Variables can be either scalar or vector. A scalar variable is a variable that contains a single value, and are created in the following form:

```
[5] a <- c(2, 3, 5)
```

Vectors can be any of the basic data types that we mentioned before. We can also apply conditions to vectors to create a logical vector, using the following conditions:

```
[6] a <- c(2, 5, 8, 3, 9)
b <- a > 3
```

We can create a vector with these values: FALSE TRUE TRUE FALSE TRUE

```
[7] b
```

FALSE TRUE TRUE FALSE TRUE

To access a particular element of the vector (vector indexing) we can use the name of the vector and the number of the element in brackets, counting from 1. The following command will produce the output corresponding to the first element in the vector a we defined:

```
[8] a[1]
```

2

We can also do logical indexing, apply a condition so that only the elements of the vector that meet the condition are shown:

```
[9] a[a>3]
```

5 8 9

Matrices

In R, we can create a matrix from a vector. Matrices are two-dimensional data structures. The way we create a matrix is by specifying the values it will contain, the number of rows and number of columns and if we are going to fill the matrix by column or by row. The following command creates a vector with values from 1 to 9.

```
[10] a <- c(1:9)
```

```
[11] A <- matrix(a, nrow=3, ncol=3, byrow=TRUE)
```

A

1	2	3
4	5	6
7	8	9

To access some element of a matrix we have to specify the row and column values in brackets after the matrix name. The following command, for example, gives the value 6, which is in the second row and the third column.

```
[12] A[2,3]
```

We can access the values of a whole column at once, by leaving an empty space where the row number would be, and vice versa. `A[2,]` gives the values of the second row. Matrices can also be created using conditions for other matrices, and can also be indexed logically. Indexing a matrix logically gives as a result a vector with the values of the matrix for which the condition is true.

```
[13] # Create a vector with values from 1 to 25
a <- c(1:25)

# Create a matrix from vector a, with 5 rows and 5 columns, and fill by row
A <- matrix(a, nrow=5, ncol=5, byrow=TRUE)

# Create a logical matrix from a condition.
C <- A > 12
```

```
[14] # Show the logical matrix: the values are true when the
# corresponding elements of matrix A meet the condition "A>12".
```

C

FALSE	FALSE	FALSE	FALSE	FALSE
FALSE	FALSE	FALSE	FALSE	FALSE
FALSE	FALSE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE
TRUE	TRUE	TRUE	TRUE	TRUE

```
[15] # Indexing matrix A using matrix C, this gives us the values
# of A for which the corresponding values of C are TRUE.
# That is, the values for which "A>12", in the form of a vector.
```

`A[C]`

```
16 21 17 22 13 18 23 14 19 24 15 20 25
```

Arrays

Arrays are data structures similar to matrices, but they can have more than two dimensions. You can create them from a vector and specify the number of dimensions of the array.

```
[16] # NOTE: For some reason, arrays don't display correctly in Jupyter not
# See below for correct output.

# Create a vector with values from 1 to 27
a <- c(1:27)

# Create an array from vector a, that has 3 matrices of dimension 3x3.
A <- array(a, dim=c(3,3,3))

# Show the array.
A
```

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
17	18	19	20	21	22	23	24	25	26	27					

This should be the right output for the array (which works in R outside Jupyter):

, , 1

	[,1]	[,2]	[,3]
[1,]	1	4	7
[2,]	2	5	8
[3,]	3	6	9

, , 2

	[,1]	[,2]	[,3]
[1,]	10	13	16
[2,]	11	14	17
[3,]	12	15	18

, , 3

	[,1]	[,2]	[,3]
[1,]	19	22	25
[2,]	20	23	26
[3,]	21	24	27

Lists

Lists are like vectors, but they can contain different data types at the same time, and they can also contain vectors.

```
[17] l <- list(c(1,2,3), 'a', 1, 1+5i)
```

```
l
```

```
1.      1    2    3
2. 'a'
3. 1
4. 1+5i
```

Data Frames

Data frames are like lists of vectors that have the same length. They are used to store data in a tabular form. You can create data frames using the following:

```
[18] data <- data.frame(
  Name = c("James", "David"),
  Gender = c('M', 'M'),
  Age = c(20, 23)
)
```

```
data
```

Name	Gender	Age
James	M	20
David	M	23

If you want to access a particular column of data of a data frame, you use the name of the data frame followed by a \$ and the name of the column. In the previous example, if we want to access the names, we can write the following, which will return a vector with the names written on it:

```
[19] data$Name
```

```
James David
```

If we want, on the other hand, access only a particular row, we can use the following, which will give the first row of the data frame:

```
[20] data[1,]
```

Name	Gender	Age
James	M	20

You can also import previously created data from other files, like CSV files, as we will see below.

Working with CMS Open Data

Importing Data from CSV files

You can import files in the CSV (comma-separated values) format to R, and analyse the data contained in them. In this tutorial we are going to use CMS Open Data as an example.

The data files for this example can be found here: <http://opendata.cern.ch/record/545>

To import the data files in R we use the following command:

```
[21] jpsi <- read.csv("http://opendata.cern.ch/record/545/files/Jpsimumu.csv")
```

Here we are saving the data contained in the `Jpsimumu.csv` file in a variable called `data`. The data files contain thousands of entries, so instead of showing the whole file, we can take a peek using the following function, which will show the first 6 rows of the data frame:

```
[22] head(jpsi)
```

Run	Event	type1	E1	px1	py1	pz1	pT
165617	75206813	G	10.1623	0.4763	-8.5164	5.5231	8.6

Run	Event	type1	E1	px1	py1	pz1	pt
165617	75678475	G	15.8799	15.0618	-1.6658	-4.7464	15.937
165617	74428554	G	21.8279	-6.2214	11.0845	17.7447	23.422
165617	75193169	G	19.4923	2.7612	-5.5769	-18.4719	19.675
165617	74832715	G	8.0972	4.6127	-1.8389	6.3949	9.425
165617	74981507	G	30.5862	15.5218	5.1293	-25.8509	35.367

Calculating the Invariant Mass

We are going to use the data stored in `jpsi`, which came from `Jpsimumu.csv`. The CSV file contains events from the CMS detector where two muons were detected.

You can see from `head(jpsi)` that there are values shown for the energy (E), the momentum (px , py , pz), the pseudo-rapidity (η or η) and the angle (ϕ or ϕ). The value for the mass is not shown, but we can calculate it from the values of the energy and the momentum.

The invariant mass can be calculated with the following equation:

$$M = \sqrt{(\sum E)^2 - \|\sum p\|^2}$$

where M is the invariant mass, $\sum E$ is the total energy, and $\sum p$ is the total momentum. To calculate the invariant mass in the code, we are going to use the values of px , py and pz and the values of the energy for the two particles. First, we are going to calculate the vector sum of the momentum of particles 1 and 2. To calculate the vector sum, we have to individually sum each component of the vectors:

```
[23] pxt <- jpsi$px1+jpsi$px2
     pyt <- jpsi$py1+jpsi$py2
     pzt <- jpsi$pz1+jpsi$pz2
```

Then we calculate the magnitude of the vector, with the following:

```
[24] # Here we use the function sqrt() to obtain the square root of
     # the sum of all momentum components squared, and save it to the
     # variable ptotal
```



```
ptotal <- sqrt(pxt^2+pyt^2+pzt^2)
```

We can also define a function to do this calculation. You can define your own functions using the following syntax:

```
myfunctionname = function(arg1, arg2...)  
{  
  statements  
  return(a)  
}
```

To define a function to calculate the magnitude of a vector, we need to take three arguments and return one value. We can call this function `magnitude()` and define it like this:

```
[25] magnitude = function(x, y, z) {  
  m = sqrt(x^2 + y^2 + z^2)  
  return(m)  
}
```

Now to do the calculation using the function we simply pass the arguments, which are the three components of the momentum in the different spatial directions:

```
[26] ptotal <- magnitude(pxt, pyt, pzt)
```

We can also define a function to do the calculation of the invariant mass of a system consisting of two particles:

```
[27] invmass = function(px1, px2, py1, py2, pz1, pz2, E1, E2){  
  px = px1+px2  
  py = py1+py2  
  pz = pz1+pz2  
  E = E1+E2  
  ptotal = magnitude(px, py, pz)  
  mass = sqrt(E^2 - ptotal^2)  
  return(mass)  
}  
  
# Here we first indicate the name of the function and its arguments,  
# then we obtain the sum of p1 and p2 for each individual component,  
# and we do the sum of the energies, then we obtain the magnitude  
# of the total momentum and finally we calculate the invariant mass  
# using the previously stated equation.
```

We can now calculate the invariant mass value using the function we have just defined and view the first values using `head()`:

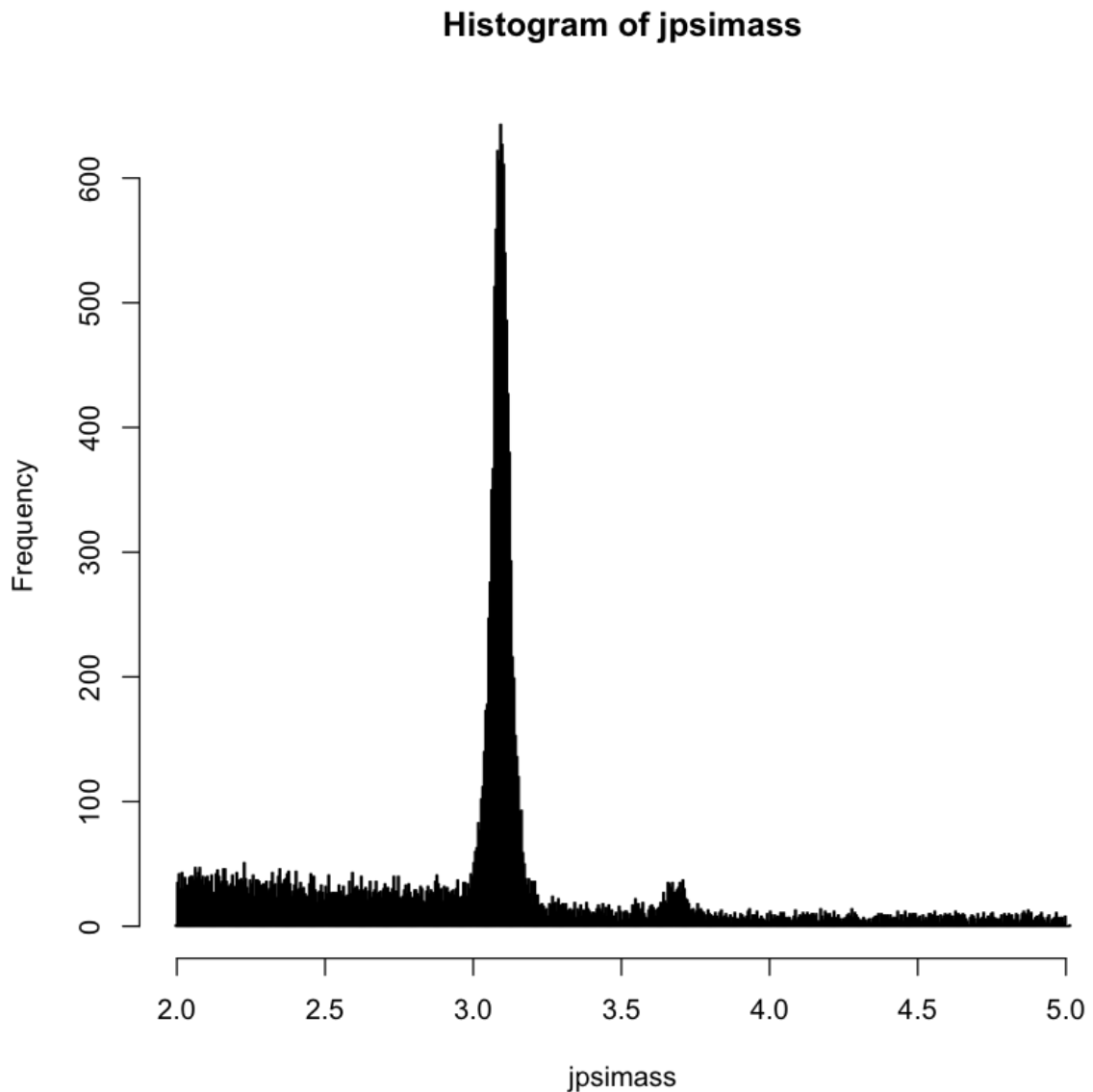
```
[28] jpsimass = invmass(jpsi$px1, jpsi$px2, jpsi$py1, jpsi$py2,  
                    jpsi$pz1, jpsi$pz2, jpsi$E1, jpsi$E2)  
  
head(jpsimass)
```

```
3.11283691027971 4.11683270974177 3.10196884252564  
2.33293443971328 4.56322815011479 3.0724626083974
```

Plotting a Histogram

Now that we have the values for the invariant mass, we can plot a histogram to see all the data. We can do this in R with the function `hist()`:

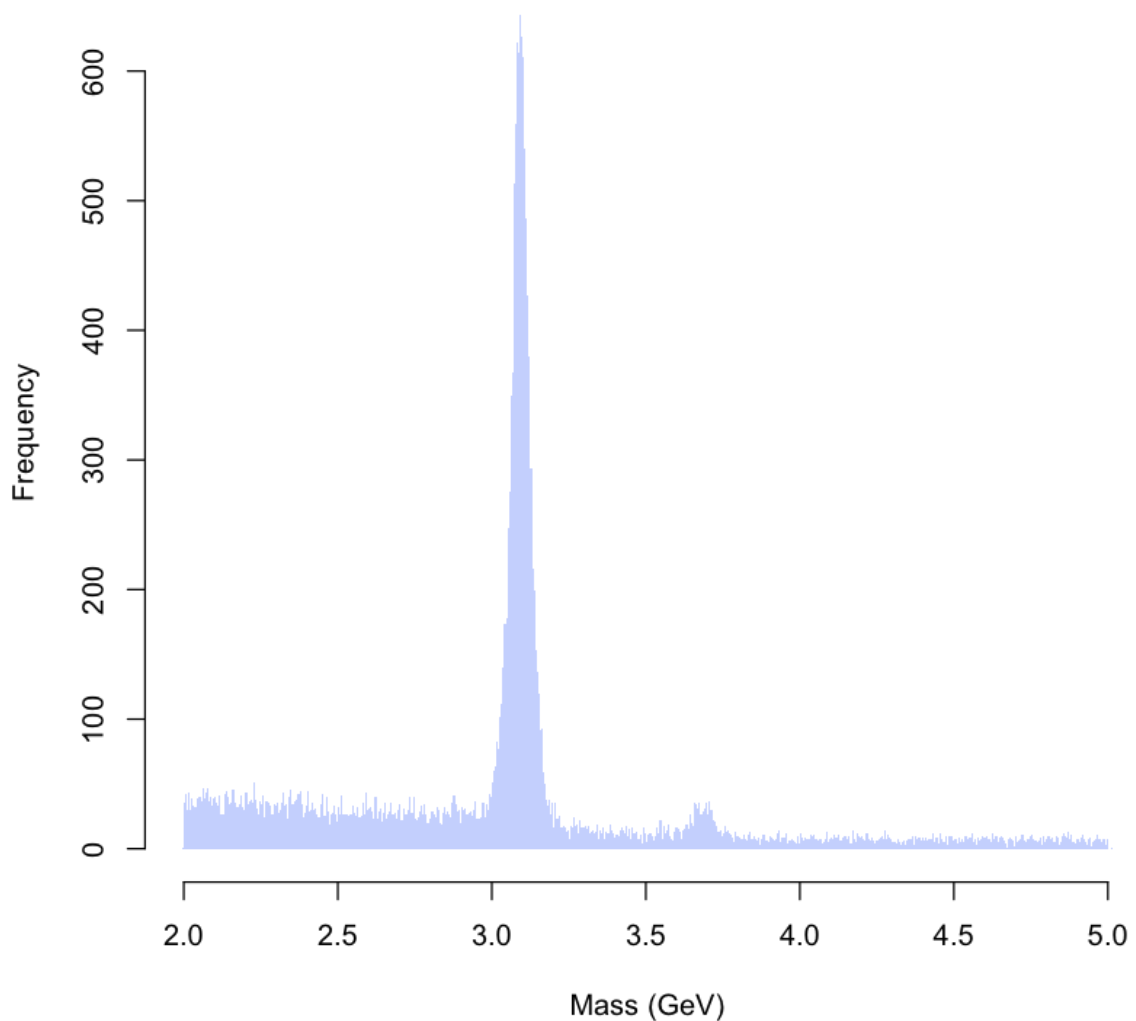
```
[29] hist(jpsimass, breaks = 500)  
  
# Here, "breaks" indicates the number of bins. You can experiment  
# changing this number to see how you can better visualise the data.
```



We can change the colors and title of the histogram with different commands. Try the following command:

```
[30] hist(jpsimass, breaks = 500, xlab = "Mass (GeV)", lty="blank",  
      col=rgb(0,0,1,1/4), main="Invariant Mass Histogram")  
  
# Here 'xlab' is the title for the x axis, lty="blank" indicates  
# that there's no border for the bins in the histogram, 'main' is  
# the title for the histogram and 'col' indicates the colour (the  
# numbers are the values for the red, green, blue and alpha channels).
```

Invariant Mass Histogram



Now in the histogram you can clearly see a very big peak around 3.1 GeV and a smaller peak around 3.7 GeV. These two peaks correspond to the mass of two particles that have di-muon decay (decay into a muon and an anti-muon). If we look in the [Particle Data Group database](#), we find that these two particles are mesons: the $J/\psi(1S)$ meson and the $\psi(2S)$ meson, respectively (although the name of the file -- `Jpsimumu.csv` -- should've given you a clue).

You can explore the rest of the files and try to identify the particles that appear in each histogram. To identify the particles, you can take into account the mass and the decay mode and search in the Particle Data Group database.