# R-bloggers

R news and tutorials contributed by hundreds of R bloggers

- [Home](#)
- [About](#)
- [RSS](#)
- [add your blog!](#)
- [Learn R](#)
- [R jobs���](#)
- [Contact us](#)

## Welcome!

Follow @rbloggers        84.5K

Here you will find daily **news and tutorials about R**, contributed by hundreds of bloggers. There are many ways to **follow us -**

By e-mail:

Your e-mail here

Subscribe

52665 readers
BY FEEDBURNER

On Facebook:

R blogg…
79K likes

Like Page

Be the first of your friends to like this

**If you are an R blogger yourself** you are invited to add your own R content feed to this site (**Non-English** R bloggers should add themselves- here)

## 🔲 Jobs for R-users

- Data Analytics Auditor, Future of Audit Lead
- Data Analytics Auditor, Future of Audit Lead @ London or Newcastle
- Senior Scientist, Translational Informatics @ Vancouver, BC

Never miss an update!
**Subscribe to R-bloggers** to receive e-mails with the latest R posts.
(You will not see this message again.)

Your E-mail…

Click here to close (This popup will not appear again)

-

# Recent Posts

-

# Other sites

-

*Never miss an update!*
**Subscribe to R-bloggers** to receive
e-mails with the latest R posts.
(You will not see this message again.)

Your E-mail…

*Click here to close (This popup will not appear again)*

# Computing and visualizing PCA in R

November 28, 2013
By [thiagogm](#)

[This article was first published on **Thiago G. Martins » R**, and kindly contributed to R-bloggers]. (You can report issue about the content on this page here)

Want to share your content on R-bloggers? click here if you have a blog, or here if you don't.

     f   Share               Tweet

Following [my introduction to PCA](#), I will demonstrate how to visualize PCA in R. There are many packages and functions apply PCA in R. In this post I will use the function `prcomp` stats package. I will also show how to visualize PCA in R R graphics. However, my favorite visualization function f `ggbiplot`, which is implemented by [Vince Q. Vu](#) and av [github](#). Please, let me know if you have better ways to visualize PCA in R.

### Computing the Principal Components (PC)

I will use the classical `iris` dataset for the demonstration. The data contain four continuous variables which corresponds to physical measures of flowers and a categorical variable describing the flowers' species.

```
1  # Load data
2  data(iris)
3  head(iris, 3)
4
5    Sepal.Length Sepal.Width Petal.Length P
6  1          5.1         3.5          1.4
7  2          4.9         3.0          1.4
8  3          4.7         3.2          1.3
```

We will apply PCA to the four continuous variables and use the categorical variable to visualize the PCs later. Notice that in the following code we apply a log transformation to the continuous variables as suggested by [1] and set `center` and `scale.` equal to `TRUE` in the call to `prcomp` to standardize the variables prior to the application of PCA:

```
1  # log transform
2  log.ir <- log(iris[, 1:4])
3  ir.species <- iris[, 5]
4
5  # apply PCA - scale. = TRUE is highly
6  # advisable, but default is FALSE.
7  ir.pca <- prcomp(log.ir,
8                   center = TRUE,
9                   scale. = TRUE)
```

Since skewness and the magnitude of the variables influence the resulting PCs, it is good practice to apply skewness transformation, center and scale the variables prior to the application of PCA. In the example above, we applied a log transformation to the variables but we could have been more general and applied a Box and Cox transformation [2]. See at the end of this post how to perform all those transformations and then apply PCA with only one call to the `preProcess` function of the `caret` package.

### Analyzing the results

The `prcomp` function returns an object of class `prcomp`, which have

deviation of each of the four PCs, and their rotation (or loadings), which
are the coefficients of the linear combinations of the continuous
variables.

```
1   # print method
2   print(ir.pca)
3
4   Standard deviations:
5   [1] 1.7124583 0.9523797 0.3647029 0.1656
6
7   Rotation:
8                      PC1          PC2
9   Sepal.Length   0.5038236 -0.45499872   0.1
10  Sepal.Width   -0.3023682 -0.88914419  -0.3
11  Petal.Length   0.5767881 -0.03378802  -0.2
12  Petal.Width    0.5674952 -0.035450
```

The plot method returns a plot of the variances (y-axis) asso
the PCs (x-axis). The Figure below is useful to decide how m
retain for further analysis. In this simple case with only 4 PC
a hard task and we can see that the first two PCs explain r
variability in the data.

```
1   # plot method
2   plot(ir.pca, type = "l")
```
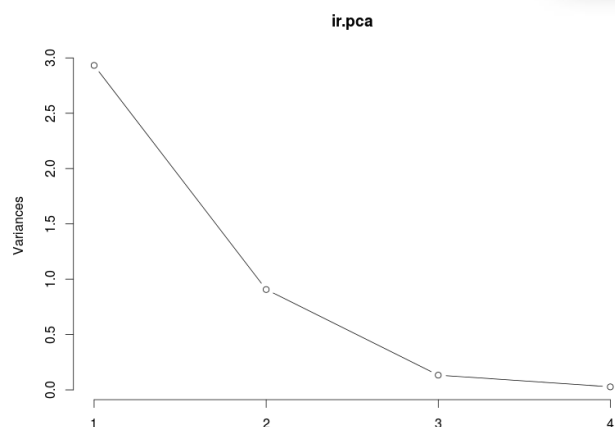


The summary method describe the importance of the PCs. The first row
describe again the standard deviation associated with each PC. The
second row shows the proportion of the variance in the data explained
by each component while the third row describe the cumulative
proportion of explained variance. We can see there that the first two PCs
accounts for more than $95\%$ of the variance of the data.

```
1   # summary method
2   summary(ir.pca)
3
4   Importance of components:
5                              PC1    PC2
6   Standard deviation      1.7125 0.9524 0.36
7   Proportion of Variance 0.7331 0.2268 0.03
8   Cumulative Proportion  0.7331 0.9599 0.99
```
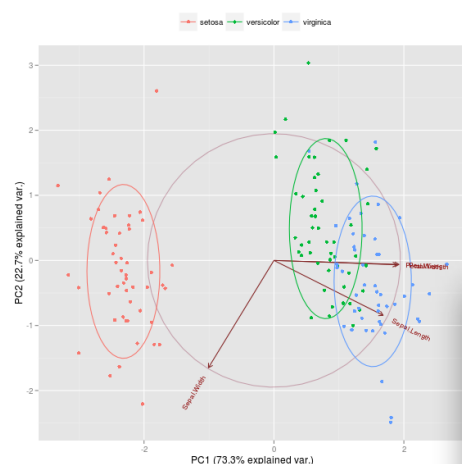
We can use the predict function if we observe new data and want to
predict their PCs values. Just for illustration pretend the last two rows of
the iris data has just arrived and we want to see what is their PCs
values:

```
1   # Predict PCs
2   predict(ir.pca,
3           newdata=tail(log.ir, 2))
4
5              PC1          PC2          PC3
6   149 1.0809930 -1.01155751 -0.7082289 -0.0
7   150 0.9712116 -0.06158655 -0.5008674 -0.1
```

The Figure below is a biplot generated by the function `ggbiplot` of the `ggbiplot` package available on [github](#).

The code to generate this Figure is given by

```
1  library(devtools)
2  install_github("ggbiplot", "vqv")
3
4  library(ggbiplot)
5  g <- ggbiplot(ir.pca, obs.scale = 1, var
6                groups = ir.species, ellip
7                circle = TRUE)
8  g <- g + scale_color_discrete(name = ''
9  g <- g + theme(legend.direction = 'hori:
10               legend.position = 'top')
11 print(g)
```

It projects the data on the first two PCs. Other PCs can be chosen through the argument `choices` of the function. It colors each point according to the flowers' species and draws a Normal contour line with `ellipse.prob` probability (default to ) for each group. More info about `ggbiplot` can be obtained by the usual `?ggbiplot`. I think you will agree that the plot produced by `ggbiplot` is much better than the one produced by `biplot(ir.pca)` (Figure below).

I also like to plot each variables coefficients inside a unit circle to get insight on a possible interpretation for PCs. Figure 4 was generated by [this code available on gist](#).

**PCA on caret package**

As I mentioned before, it is possible to first apply a Box-Cox transformation to correct for skewness, center and scale each variable and then apply PCA in one call to the `preProcess` function of the `caret` package.

```
1  require(caret)
2  trans = preProcess(iris[,1:4],
3                method=c("BoxCox", "c
4                         "scale", "pc:
5  PC = predict(trans, iris[,1:4])
```

By default, the function keeps only the PCs that are necessary to explain at least 95% of the variability in the data, but this can be changed through the argument `thresh`.

```
1  # Retained PCs
2  head(PC, 3)
3
4         PC1        PC2
5  1 -2.303540 -0.4748260
6  ? ? 1?1?10  0 ?4??003
```

```
 9   # Loadings
10   trans$rotation
11
12                    PC1          PC2
13   Sepal.Length  0.5202351 -0.38632246
14   Sepal.Width  -0.2720448 -0.92031253
15   Petal.Length  0.5775402 -0.04885509
16   Petal.Width   0.5672693 -0.03732262
```

See [Unsupervised data pre-processing for predictive modeling](#) for an introduction of the `preProcess` function.

**References:**

[1] Venables, W. N., Brian D. R. Modern applied statistics with S-PLUS. Springer-verlag. (Section 11.1)
[2] Box, G. and Cox, D. (1964). An analysis of transformatio[...]
of the Royal Statistical Society. Series B (Methodological) 21[...]

f  Share        🐦 Tweet

To **leave a comment** for the author, please follow the link and comment on their blog: **Thiago G. Martins » R**.

R-bloggers.com offers **daily e-mail updates** about R news and tutorials about learning R and many other topics. Click here if you're looking to post or find an R/data-science job.

Want to share your content on R-bloggers? click here if you have a blog, or here if you don't.

---

If you got this far, why not **subscribe for updates** from the site? Choose your flavor: e-mail, twitter, RSS, or facebook...

Like 48    Share      Tweet     Share

Comments are closed.

# Search R-bloggers

Search.. | Go

# Most visited articles of the week

1. Free Springer Books during COVID19
2. 5 Ways to Subset a Data Frame in R
3. How to write the first for loop in R
4. R – Sorting a data frame by the contents of a column
5. Date Formats in R
6. How to schedule R scripts
7. Installing R packages
8. In-depth introduction to machine learning in 15 hours of expert videos
9. Intro to {polite} Web Scraping of Soccer Data with R!

# Sponsors