# Analyze the internet data of www.datadb.com

*Business Analytic Foundation with R Tools- Question*

## Abstract

The web analytics team of www.datadb.com is interested to understand the web activities of the site, which are the sources used to access the website. They have a database that states the keywords of time in page, source group, bounces, exits, unique page views, and visits.

## Presented by: Ankita Agarwal

# Problem Statement:

The web analytics team of www.datadb.com is interested to understand the web activities of the site, which are the sources used to access the website. They have a database that states the keywords of time in page, source group, bounces, exits, unique page views, and visits.

# Detailed description of the given dataset:

**Bounces:** It represents the percentage of visitors who enter the site and "bounce" (leave the site) rather than continuing to view other pages within the same site.

**Exits:** It represents the percentage of visitors to a site who actively click away to a different site from a specific page, after possibly having visited any other page on the site.

**Continent:** It shows the continent from which the site has been accessed.

**Source group:** It shows how the visitor has accessed the site.

**Time on page:** It shows how long the user has spent on that particular page of the website.

**Unique page view:** It represents the number of sessions during which that page was viewed one or more times.

**Visits:** A visit counts all visitors, no matter how many times the same visitor may have been to your site.
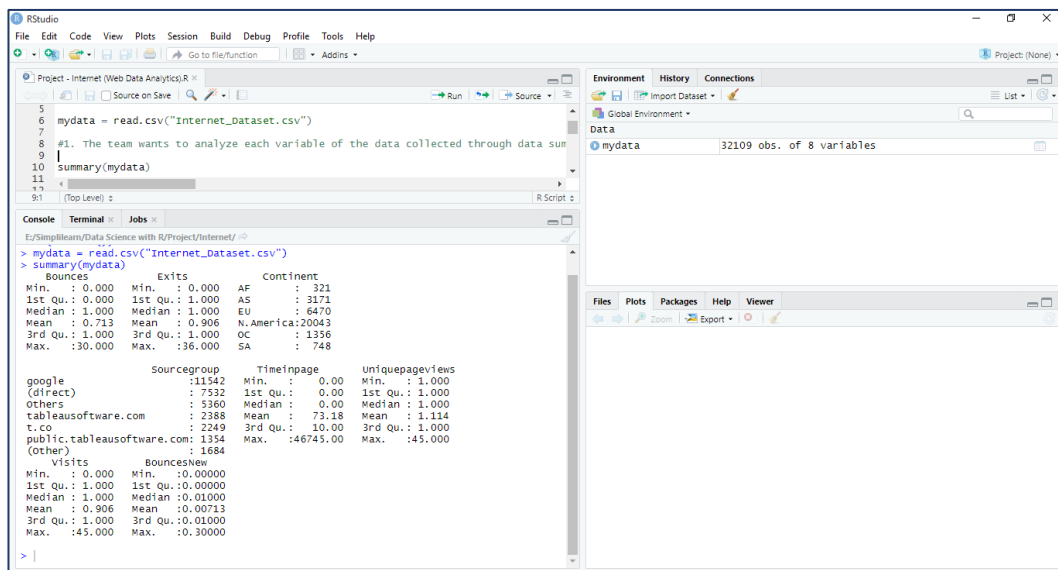
# To Analyze:

1. The team wants to analyze each variable of the data collected through data summarization to get a basic understanding of the dataset and to prepare for further analysis.

2. As mentioned earlier, a unique page view represents the number of sessions during which that page was viewed one or more times. A visit counts all instances, no matter how many times the same visitor may have been to your site. So, the team needs to know whether the unique page view value depends on visits.

3. Find out the probable factors from the dataset, which could affect the exits. Exit Page Analysis is usually required to get an idea about why a user leaves the website for a session and moves on to another one. Please keep in mind that exits should not be confused with bounces.

4. Every site wants to increase the time on page for a visitor. This increases the chances of the visitor understanding the site content better and hence there are more chances of a transaction taking place. Find the variables which possibly have an effect on the time on page.

5. A high bounce rate is a cause of alarm for websites which depend on visitor engagement. Help the team in determining the factors that are impacting the bounce.

# Analysis and Interpretations:

1. **The team wants to analyze each variable of the data collected through data summarization to get a basic understanding of the dataset and to prepare for further analysis.**

To get a basic insight into the data, we do a 5-point summary analysis for the data set.



**Interpretation:**

By looking at the 5-point summary of the data, we can summarize that percentage of visitors that bounce or exit through the website ranges from 0 to 30 and 36 respectively, however, maximum visitors only do it once as shown by the median of the variables. Max number of visitors belong to North America – 20043 or accessed it using google – 11542. The average time spent on the page comes to be 73.2 seconds.
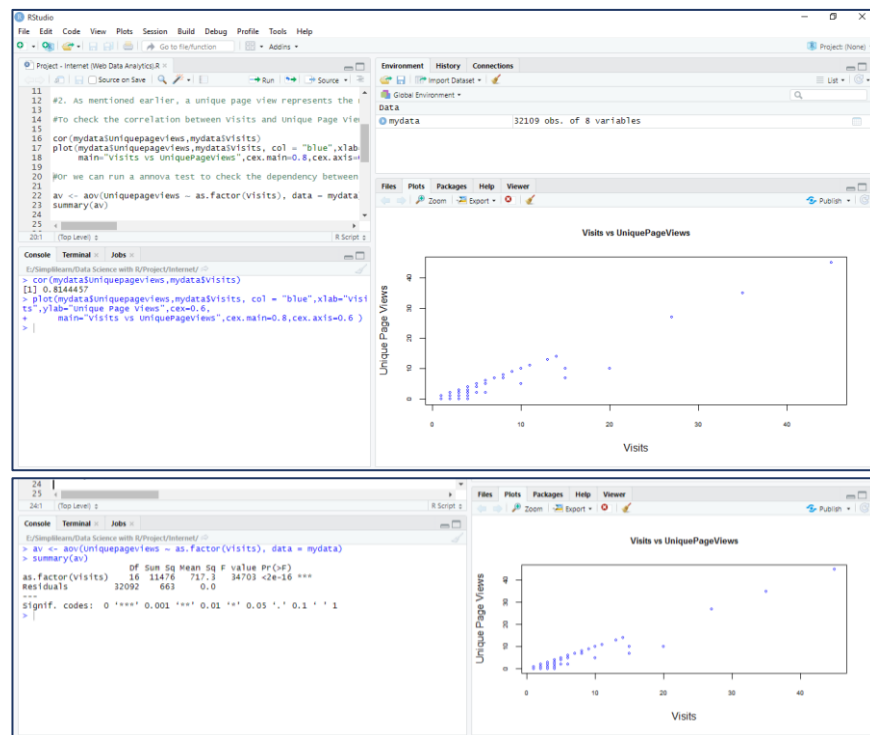
## 2. As mentioned earlier, a unique page view represents the number of sessions during which that page was viewed one or more times. A visit counts all instances, no matter how many times the same visitor may have been to your site. So, the team needs to know whether the unique page view value depends on visits.

In order to check the dependencies of the variables to each other (unique page view depends on visits) we can either create check the correlation between the variables or do an ANOVA test to understand the dependency between the variables.

Defining Hypothesis:

Ho: The unique page view value does not depend on visits

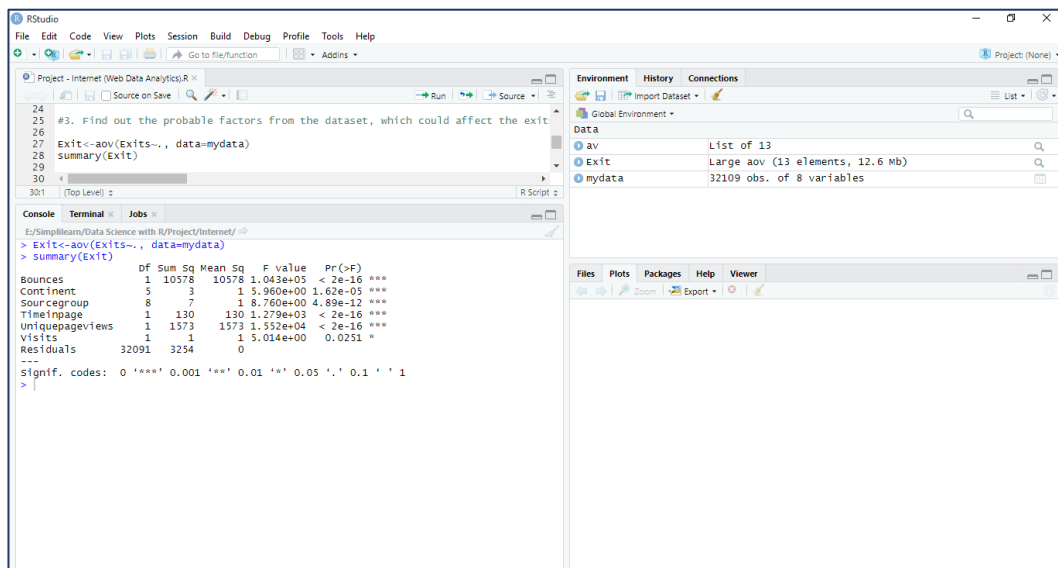Ha: The unique page view value depends on visits



**Interpretation:**

As seen in the above plot, there is a positive linear correlation between the two variables Unique page views and visits.

Also, we know, when p-value < alpha, p-value is less than alpha; we reject the null hypothesis. We take alpha value as 0.05 at 95% confidence level. Here, p-value significantly less than the speculated 0.05, concluding that the variable visits have a very high impact on Unique page views.

## 3. Find out the probable factors from the dataset, which could affect the exits. Exit Page Analysis is usually required to get an idea about why a user leaves the website for a session and moves on to another one. Please keep in mind that exits should not be confused with bounces.

To check the dependencies of other variables on Exits we do an ANOVA test to understand the dependencies between the variables.
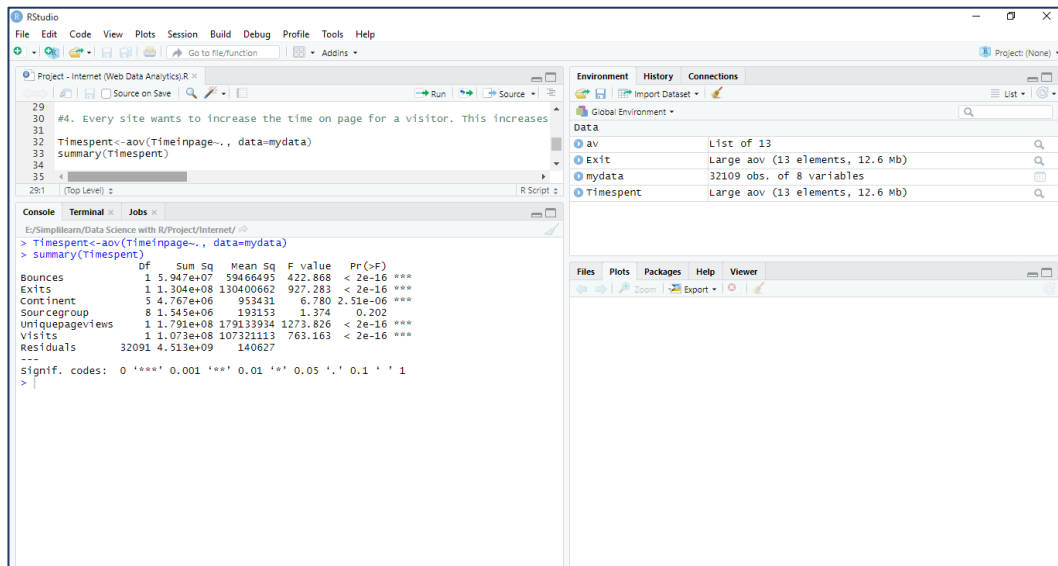


**Interpretation:**

As we know, when p-value is less than 0.05, which we assume while creating regression model, we consider that variable to be significant in the model.

Here, we can see that Bounces, Unique page views and Time in page have significantly low p-values, implying that these variables are extremely significant to exit. Although Continent and Source group have low p-values, they are less significant as compared to bounces, time in page and unique page views.

Also, number of visits do not have as much significance as other variables.

.

**4. Every site wants to increase the time on page for a visitor. This increases the chances of the visitor understanding the site content better and hence there are more chances of a transaction taking place. Find the variables which possibly have an effect on the time on page.**

To check the dependencies of other variables on Time spent in page we do an ANOVA test to understand the dependencies between the variables.
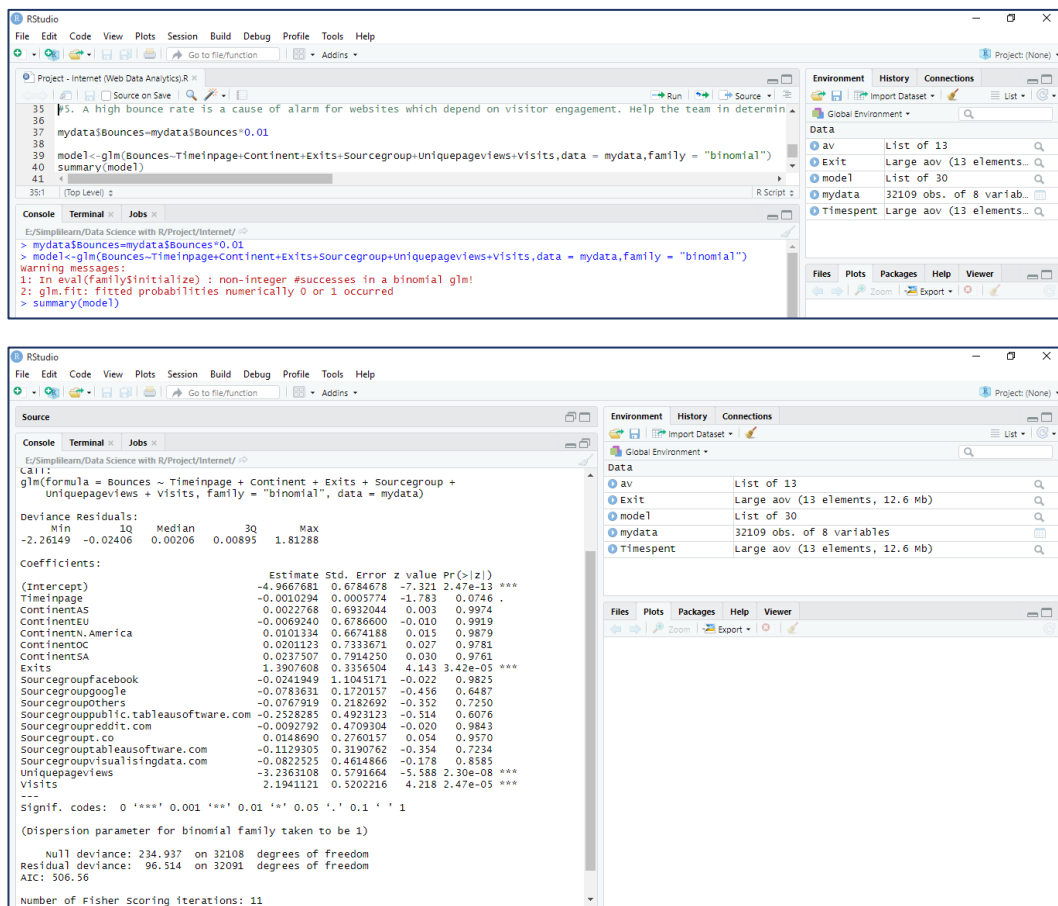


**Interpretation:**

As we know, when p-value is less than 0.05, which we assume while creating regression model, we consider that variable to be significant in the model.

Here, we can see that Bounces, Exits, Unique page views and Visits have significantly low p-values, implying that these variables are extremely significant to Time in page. Although Continent has a low p-value, it is less significant.

As noticed only source group with high p-value does not have much significance to Time spent in the page.

## 5. A high bounce rate is a cause of alarm for websites which depend on visitor engagement. Help the team in determining the factors that are impacting the bounce.

To understand factors impacting the bounce we create a logistic regression model. The data variables for bounce should be between 0 and 1 as logistic regression shows probability which is always between 0-1. Hence, we multiply bounce with 0.01 to create a regression model.





**Interpretation:**

As we know, when p-value is less than 0.05, which we assume while creating regression model, we consider that variable to be significant in the model.

Here, we can see Exists, Unique page views and Visits have lesser p-value, implying that these variables are extremely significant to bounce.

# Programming Codes:

#Reading Internet Data

mydata = read.csv("Internet_Dataset.csv")

#1.     The team wants to analyze each variable of the data collected through data summarization to get a basic understanding of the dataset and to prepare for further analysis.

summary(mydata)

#2.     As mentioned earlier, a unique page view represents the number of sessions during which that page was viewed one or more times. A visit counts all instances, no matter how many times the same visitor may have been to your site. So, the team needs to know whether the unique page view value depends on visits.

#To check the correlation between Visits and Unique Page Views

```
cor(mydata$Uniquepageviews,mydata$Visits)
plot(mydata$Uniquepageviews,mydata$Visits, col = "blue",xlab="Visits",ylab="Unique Page Views",cex=0.6,
    main="Visits vs UniquePageViews",cex.main=0.8,cex.axis=0.6 )
```

#Or we can run a annova test to check the dependency between Visits and Unique Page Views

```
av <- aov(Uniquepageviews ~ as.factor(Visits), data = mydata)
summary(av)
```

#3.     Find out the probable factors from the dataset, which could affect the exits. Exit Page Analysis is usually required to get an idea about why a user leaves the website for a session and moves on to another one. Please keep in mind that exits should not be confused with bounces.

```
Exit<-aov(Exits~., data=mydata)
summary(Exit)
```

#4.     Every site wants to increase the time on page for a visitor. This increases the chances of the visitor understanding the site content better and hence there are more chances of a transaction taking place. Find the variables which possibly have an effect on the time on page.

```
Timespent<-aov(Timeinpage~., data=mydata)
summary(Timespent)
```

#5.     A high bounce rate is a cause of alarm for websites which depend on visitor engagement. Help the team in determining the factors that are impacting the bounce.

```
mydata$Bounces=mydata$Bounces*0.01

model<-glm(Bounces~Timeinpage+Continent+Exits+Sourcegroup+Uniquepageviews+Visits,data =
mydata,family = "binomial")
summary(model)
```

-------------------------------------------------------------------The End-------------------------------------------------------------------