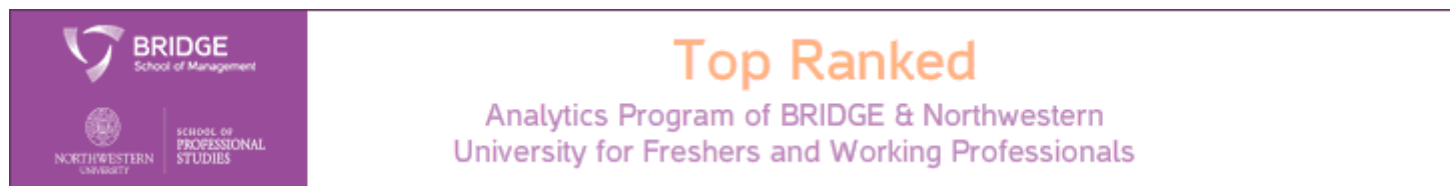


f (<https://www.facebook.com/AnalyticsVidhya>)t (<https://twitter.com/analyticsvidhya>)g+ (<https://plus.google.com/+Analyticsvidhya/posts>)in (<https://www.linkedin.com/groups/Analytics-Vidhya-Learn-everything-about-5057165>)(<https://datahack.analyticsvidhya.com/contest/skilltest-statistics/>)Home (<https://www.analyticsvidhya.com/>) > R (<https://www.analyticsvidhya.com/blog/category/r/>) > Tutorial on 5 Powerful R Packages used for imputing missing values (htt..

Tutorial on 5 Powerful R Packages used for imputing missing values

R (<https://www.analyticsvidhya.com/blog/category/r/>)

SHARE f (<http://www.facebook.com/sharer.php?u=https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/&t=Tutorial%20on%205%20Powerful%20R%20Packages%20used%20for%20imputing%20missing%20values>) t (<https://twitter.com/home?status=Tutorial%20on%205%20Powerful%20R%20Packages%20used%20for%20imputing%20missing%20values+https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/>) g+ (<https://plus.google.com/share?url=https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/>) p (<http://pinterest.com/pin/create/button/?url=https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/&media=https://www.analyticsvidhya.com/wp-content/uploads/2016/03/imputing.jpg&description=Tutorial%20on%205%20Powerful%20R%20Packages%20used%20for%20imputing%20missing%20values>)



(http://admissions.bridgesom.com/pba-new/?utm_source=AV&utm_medium=BannerInline&utm_campaign=AVBanner20August)

Introduction

Missing values are considered to be the first obstacle in predictive modeling. Hence, it's important to master the methods to overcome them. Though, some machine learning algorithms (<https://www.analyticsvidhya.com/blog/2015/09/random-forest-algorithm-multiple-challenges/>) claim to treat them intrinsically, but who knows how good it happens inside the 'black box'.

The choice of method to impute missing values, largely influences the model's predictive ability. In most statistical analysis methods, listwise deletion is the default method used to impute missing values. But, it not as good since it leads to information loss.

Do you know R has robust packages for missing value imputations?

Yes! R Users have something to cheer about. We are endowed with some incredible R packages for missing values imputation. These packages arrive with some inbuilt functions and a simple syntax to impute missing data at once. Some packages are known best working with continuous variables and others for categorical. With this article, you can make a better decision choose the best suited package.

In this article, I've listed 5 R packages popularly known for missing value imputation. There might be more packages. But, I decided to focus on these ones. I've tried to explain the concepts in simplistic manner with practice examples in R.



List of R Packages

1. MICE
2. Amelia
3. missForest
4. Hmisc
5. mi

MICE Package

MICE (Multivariate Imputation via Chained Equations) is one of the commonly used package by R users. Creating multiple imputations as compared to a single imputation (such as mean) takes care of uncertainty in missing values.

MICE assumes that the missing data are Missing at Random (MAR), which means that the probability that a value is missing depends only on observed value and can be predicted using them. It imputes data on a variable by variable basis by specifying an imputation model per variable.

For example: Suppose we have X_1, X_2, \dots, X_k variables. If X_1 has missing values, then it will be regressed on other variables X_2 to X_k . The missing values in X_1 will be then replaced by predictive values obtained. Similarly, if X_2 has missing values, then X_1, X_3 to X_k variables will be used in prediction model as independent variables. Later, missing values will be replaced with predicted values.

By default, linear regression is used to predict continuous missing values. Logistic regression is used for categorical missing values. Once this cycle is complete, multiple data sets are generated. These data sets differ only in imputed missing values. Generally, it's considered to be a good practice to build models on these data sets separately and combining their results.

Precisely, the methods used by this package are:

1. PMM (Predictive Mean Matching) – For numeric variables
2. logreg(Logistic Regression) – For Binary Variables(with 2 levels)
3. polyreg(Bayesian polytomous regression) – For Factor Variables (≥ 2 levels)
4. Proportional odds model (ordered, ≥ 2 levels)

Let's understand it practically now.

```
> path <- "../Data/Tutorial"  
> setwd(path)
```

```
#load data
> data <- iris

#Get summary
> summary(iris)
```

Since, MICE assumes missing at random values. Let's seed missing values in our data set using `prodNA` function. You can access this function by installing `missForest` package.

```
#Generate 10% missing values at Random
> iris.mis <- prodNA(iris, noNA = 0.1)

#Check missing values introduced in the data
> summary(iris.mis)
```

I've removed categorical variable. Let's here focus on continuous values. To treat categorical variable, simply encode the levels and follow the procedure below.

```
#remove categorical variables
> iris.mis <- subset(iris.mis, select = -c(Species))
> summary(iris.mis)

#install MICE
> install.packages("mice")
> library(mice)
```

`mice` package has a function known as *md.pattern()*. It returns a tabular form of missing value present in each variable in a data set.

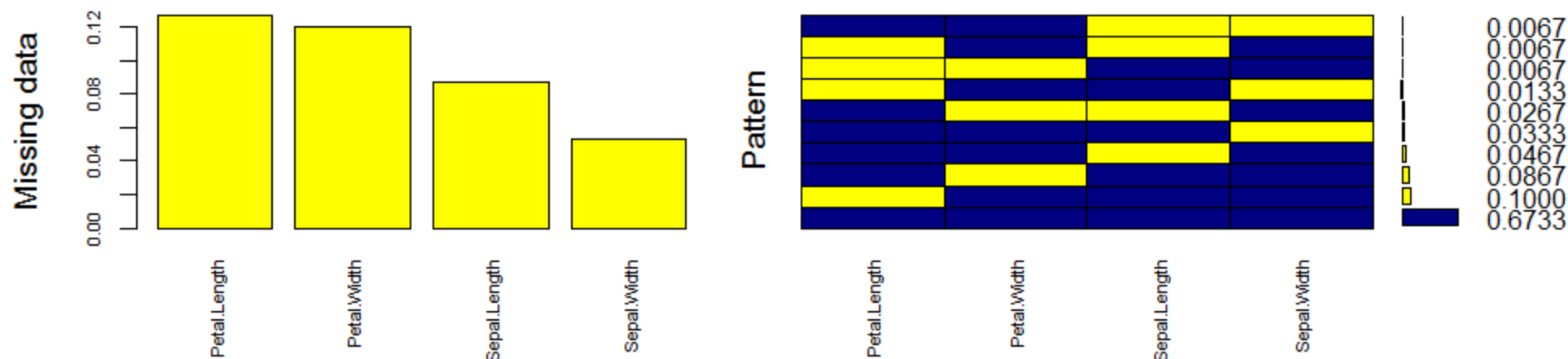
```
> md.pattern(iris.mis)
```

	Sepal.Length	Sepal.Width	Petal.Width	Petal.Length	
98	1	1	1	1	0
10	0	1	1	1	1
13	1	0	1	1	1
12	1	1	1	0	1
12	1	1	0	1	1
2	0	1	1	0	2
1	1	0	0	1	2
1	1	1	0	0	2
1	0	1	0	0	3
	13	14	15	16	58

Let's understand this table. There are 98 observations with no missing values. There are 10 observations with missing values in Sepal.Length. Similarly, there are 13 missing values with Sepal.Width and so on.

This looks ugly. Right ? We can also create a visual which represents missing values. It looks pretty cool too. Let's check it out.

```
> install.packages("VIM")
> library(VIM)
> mice_plot <- aggr(iris.mis, col=c('navyblue','yellow'),
  numbers=TRUE, sortVars=TRUE,
  labels=names(iris.mis), cex.axis=.7,
  gap=3, ylab=c("Missing data","Pattern"))
```



Let's quickly understand this. There are 67% values in the data set with no missing value. There are 10% missing values in Petal.Length, 8% missing values in Petal.Width and so on. You can also look at histogram which clearly depicts the influence of missing values in the variables.

Now, let's impute the missing values.

```
> imputed_Data <- mice(iris.mis, m=5, maxit = 50, method = 'pmm', seed = 500)
> summary(imputed_Data)
```

Multiply imputed data set

Call:

```
mice(data = iris.mis, m = 5, method = "pmm", maxit = 50, seed = 500)
```

Number of multiple imputations: 5

Missing cells per column:

```
Sepal.Length Sepal.Width Petal.Length Petal.Width
```

```
13           14           16           15
```

Imputation methods:

```
Sepal.Length Sepal.Width Petal.Length Petal.Width
```

```
"pmm"        "pmm"        "pmm"        "pmm"
```

VisitSequence:

```
Sepal.Length Sepal.Width Petal.Length Petal.Width
```

```
1           2           3           4
```

PredictorMatrix:

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
Sepal.Length	0	1	1	1
Sepal.Width	1	0	1	1
Petal.Length	1	1	0	1
Petal.Width	1	1	1	0

Random generator seed value: 500

Here is an explanation of the parameters used:

1. m – Refers to 5 imputed data sets
2. maxit – Refers to no. of iterations taken to impute missing values
3. method – Refers to method used in imputation. we used predictive mean matching.

#check imputed values



can select any using *complete()* function.



(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

```
> fit <- with(data = iris.mis, exp = lm(Sepal.Width ~ Sepal.Length + Petal.Width))
```

```
#combine results of all 5 models
```

```
> combine <- pool(fit)
```

```
> summary(combine)
```

Please note that I've used the command above just for demonstration purpose. You can replace the variable values at your end and try it.

Amelia

This package (Amelia II) is named after Amelia Earhart, the first female aviator to fly solo across the Atlantic Ocean. History says, she got mysteriously disappeared (missing) while flying over the pacific ocean in 1927, hence this package was named to solve missing value problems.



imputation (generate imputed data sets) to deal with missing data. It helps to reduce bias and increase efficiency. It is enabled with bootstrap and is faster and robust to impute many variables including cross





enabled with parallel imputation feature using multicore CPUs.

Multivariate Normal Distribution (MVN). It uses means and covariances to summarize data. (Sampling at Random)

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

It works this way. First, it takes m bootstrap samples and applies EMB algorithm to each sample. The m estimates of mean and variances will be different. Finally, the first set of estimates are used to impute first set of missing values using regression, then second set of estimates are used for second set and so on.

On comparing with MICE, MVN lags on some crucial aspects such as:

1. MICE imputes data on variable by variable basis whereas MVN uses a joint modeling approach based on multivariate normal distribution.
2. MICE is capable of handling different types of variables whereas the variables in MVN need to be normally distributed or transformed to approximate normality.
3. Also, MICE can manage imputation of variables defined on a subset of data whereas MVN cannot.

Hence, this package works best when data has multivariable normal distribution. If not, transformation is to be done to bring $\hat{\text{data}}$ close to normality.

Let's understand it practically now.

```
#install package and load library
```

```
> install.packages("Amelia")
```





ul about is classifying variables. It has 3 parameters:

er variables which you don't want to impute

```
> iris.mis <- prodNA(iris, nonA = 0, 1)
(https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)
```

#specify columns and run amelia

```
> amelia_fit <- amelia(iris.mis, m=5, parallel = "multicore", noms = "Species")
```

#access imputed outputs

```
> amelia_fit$imputations[[1]]
```

```
> amelia_fit$imputations[[2]]
```

```
> amelia_fit$imputations[[3]]
```

```
> amelia_fit$imputations[[4]]
```

```
> amelia_fit$imputations[[5]]
```

To check a particular column in a data set, use the following commands

```
> amelia_fit$imputations[[5]]$Sepal.Length
```



```
= "imputed_data_set")
```



implementation of random forest (<https://www.analyticsvidhya.com/blog/2015/09/random-forest-imputation/>). It's a non parametric imputation method applicable to various variable types. So,

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

Non-parametric method does not make explicit assumptions about functional form of f (any arbitrary function). Instead, it tries to estimate f such that it can be as close to the data points without seeming impractical.

How does it work ? In simple words, it builds a random forest model for each variable. Then it uses the model to predict missing values in the variable with the help of observed values.

It yield OOB (out of bag) imputation error estimate. Moreover, it provides high level of control on imputation process. It has options to return OOB separately (for each variable) instead of aggregating over the whole data matrix. This helps to look more closely as to how accurately the model has imputed values for each variable.

Let's understand it practically. Since bagging works well on categorical variable too, we don't need to remove them here. It very well takes care of missing value pertaining to their variable types:

```
#missForest
```

```
> install.packages("missForest")
```

```
> library(missForest)
```





parameters as default values

```
> iris.imp <- missForest(iris.mis)
(https://datahack.analyticsvidhya.com/contest/av-
```

```
casino-introduction-to-probability/)
```

```
#check imputed values
```

```
> iris.imp$ximp
```

```
#check imputation error
```

```
> iris.imp$OOBerror
```

```
NRMSE      PFC
```

```
0.14148554 0.02985075
```

NRMSE is normalized mean squared error. It is used to represent error derived from imputing continuous values. PFC (proportion of falsely classified) is used to represent error derived from imputing categorical values.

```
#comparing actual data accuracy
```

```
> iris.err <- mixError(iris.imp$ximp, iris.mis, iris)
```

```
>iris.err
```





(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

are imputed with 6% error and continuous variables are imputed with 15% error. This can be controlled by the `ntree` parameter. `mtry` refers to the number of variables being randomly sampled at each node in the forest.

Hmisc is a multiple purpose package useful for data analysis, high – level graphics, imputing missing values, advanced table making, model fitting & diagnostics (linear regression, logistic regression & cox regression) etc. Amidst, the wide range of functions contained in this package, it offers 2 powerful functions for imputing missing values. These are `impute()` and `aregImpute()`. Though, it also has `transcan()` function, but `aregImpute()` is better to use.

`impute()` function simply imputes missing value using user defined statistical method (mean, max, median). Its default is median. On the other hand, `aregImpute()` allows mean imputation using additive regression, bootstrapping, and predictive mean matching.

In bootstrapping, different bootstrap resamples are used for each of multiple imputations. Then, a flexible additive model ($\hat{\mu}$ non parametric regression method) is fitted on samples taken with replacements from original data and missing values (acts as dependent variable) are predicted using non-missing values (independent variable).

Then, it uses predictive mean matching (default) to impute missing values. Predictive mean matching works well for continuous and categorical (binary & multi-level) without the need for computing residuals and maximum likelihood fit.



is package:

being predicted.

(https://en.wikipedia.org/wiki/Scoring_algorithm) method is used for predicting categorical variables.



```
(https://datahack.analyticsvidhya.com/contest/av-
#load data
casino-introduction-to-probability/)
> data("iris")
```

```
#seed missing values ( 10% )
```

```
> iris.mis <- prodNA(iris, noNA = 0.1)
```

```
> summary(iris.mis)
```

```
# impute with mean value
```

```
> iris.mis$imputed_age <- with(iris.mis, impute(Sepal.Length, mean))
```

```
# impute with random value
```

```
> iris.mis$imputed_age2 <- with(iris.mis, impute(Sepal.Length, 'random'))
```

```
#similarly you can use min, max, median to impute missing value
```

```
#using argImpute
```

```
> impute_age <- argImpute( Sepal.Length + Sepal.Width + Petal.Length + Petal.Width +
)
```



variable type and treats them accordingly.



n: 150 p: 5 Imputations: 5 nk: 3
<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>

number of last

Sepal.Length	Sepal.width	Petal.Length	Petal.width	Species
21	12	12	14	16

	type	d.f.
Sepal.Length	s	2
Sepal.width	s	2
Petal.Length	s	2
Petal.width	s	2
Species	c	2

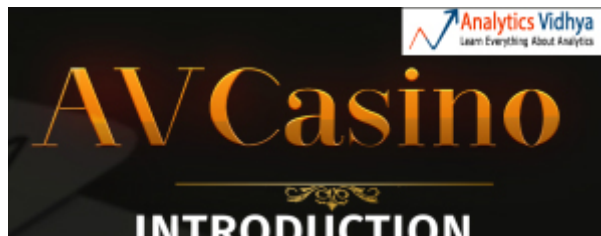
Transformation of Target Variables Forced to be Linear

R-squares for Predicting Non-Missing Values for Each Variable
 Using Last Imputations of Predictors

Sepal.Length	Sepal.width	Petal.Length	Petal.width	Species
0.865	0.670	0.984	0.958	0.988

The output shows R^2 values for predicted missing values. Higher the value, better are the values predicted. You can also check imputed values using the following command

```
#check imputed variable Sepal.Length
```





mi) package provides several features for dealing with missing values. Like other packages, it uses predictive mean matching (pmm) to approximate missing values. And, uses predictive mean matching method.

Though, I've already explained predictive mean matching (pmm) above, but if you haven't understood yet, here's a simpler version: For each observation in a variable with missing value, we find observation (from available values) with the closest predictive mean to that variable. The observed value from this "match" is then used as imputed value.

Below are some unique characteristics of this package:

1. It allows graphical diagnostics of imputation models and convergence of imputation process.
2. It uses bayesian version of regression models to handle issue of separation.
3. Imputation model specification is similar to regression output in R
4. It automatically detects irregularities in data such as high collinearity among variables.
5. Also, it adds noise to imputation process to solve the problem of additive constraints.

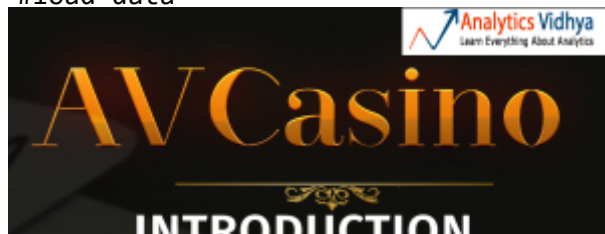
Let's understand it practically.

#install package and load library

```
> install.packages("mi")
```

```
> library(mi)
```

#load data





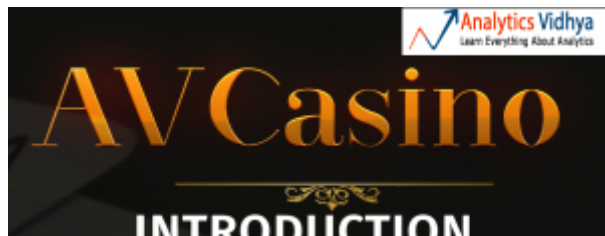
```
> mi_data <- mi(iris.mis, seed = 335)
```

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)
I've used default values of parameters namely:

1. rand.imp.method as "bootstrap"
2. n.imp (number of multiple imputations) as 3
3. n.iter (number of iterations) as 30

```
> summary(mi_data)
```

^





INTRODUCTION
to
PROBABILITY
2016, Aug 10th - 31st
ONLINE **REGISTER NOW**

u. Max.
-1.0900 -0.0410 -0.4038 -0.2237 0.1647 1.5550

<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>

`$sepal.length$observed`

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-0.93460	-0.43070	-0.05273	0.00000	0.32520	1.14400

`$sepal.width`
`$sepal.width$is_missing`
missing
FALSE TRUE
138 12

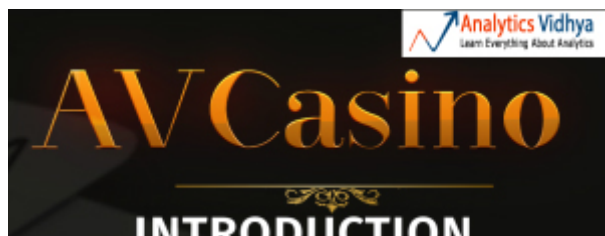
`$sepal.width$imputed`

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-0.85220	-0.23360	0.08939	0.08501	0.40860	1.30000

`$sepal.width$observed`

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-1.23600	-0.30270	-0.06934	0.00000	0.28070	1.56400

Here is a snapshot of summary output by `mi` package after imputing missing values. As shown, it uses summary statistics to define the imputed values.





...s ? I am sure many of you would be asking this! Having created this tutorial, I felt Hmisc should
 ...utation followed by missForest and MICE.

...variables types and uses bootstrap sample and predictive mean matching to impute missing
 ...eat categorical variable, just like we did while using MICE package. However, missForest can
 ...utperform MICE if the observed variables supplied contain sufficient information.

([https://datahack.analyticsvidhya.com/contest/av-](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

[casino-introduction-to-probability/](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/))

In this article, I explain using 5 different R packages for missing value imputation. Such advanced methods can help you score better
 accuracy in building predictive models.

Did you find this article useful ? Which package do you generally use to impute missing values ? Do share your experience /
 suggestions in the comments section below.

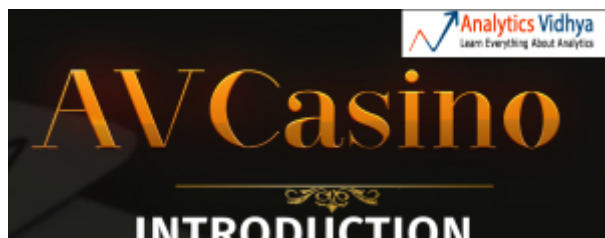
**You want to apply your analytical skills and test your potential? Then participate in our Hackathons
 (<http://datahack.analyticsvidhya.com/contest/all>) and compete with Top Data Scientists from all over
 the world.**

Share this:

 (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?share=linkedin&nb=1>) 450

 (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?share=facebook&nb=1>)

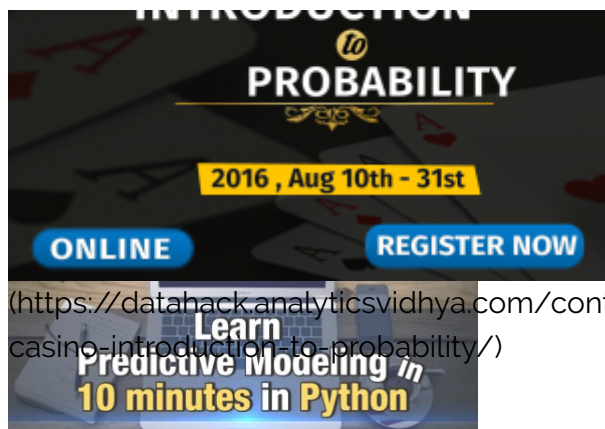
 (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?share=google-plus-1&nb=1>)



[tutorial-powerful-packages-imputing-missing-values/?share=twitter&nb=1](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?share=twitter&nb=1))

[tutorial-powerful-packages-imputing-missing-values/?share=pocket&nb=1](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?share=pocket&nb=1))

[tutorial-powerful-packages-imputing-missing-values/?share=reddit&nb=1](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?share=reddit&nb=1))



(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

(<https://www.analyticsvidhya.com/blog/2015/09/build-predictive-model-10-minutes-python/>)
Build a Predictive Model in 10 Minutes (using Python)

(<https://www.analyticsvidhya.com/blog/2015/09/build-predictive-model-10-minutes-python/>)

In "Business Analytics"



(<https://www.analyticsvidhya.com/blog/2016/03/select-important-variables-boruta-package/>)

How to perform feature selection (i.e. pick important variables) using Boruta Package in R ?

(<https://www.analyticsvidhya.com/blog/2016/03/select-important-variables-boruta-package/>)

In "R"



(<https://www.analyticsvidhya.com/blog/2015/02/7-steps-data-exploration-preparation-building-model-part-2/>)

7 Steps of Data Exploration & Preparation - Part 2

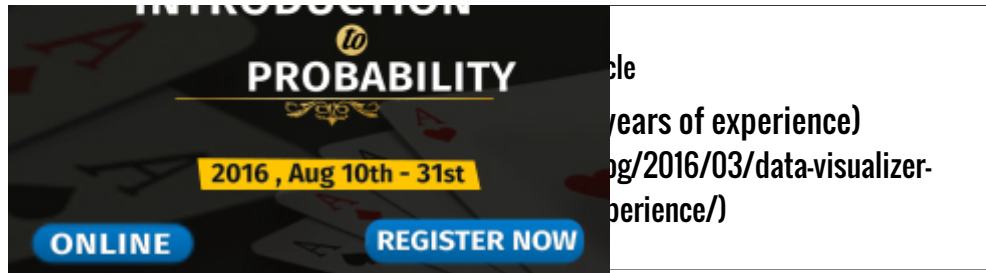
(<https://www.analyticsvidhya.com/blog/2015/02/7-steps-data-exploration-preparation-building-model-part-2/>)

In "Business Analytics"

TAGS: AMELIA PACKAGE (<https://www.analyticsvidhya.com/blog/tag/amelia-package/>), BOOTSTRAP SAMPLING (<https://www.analyticsvidhya.com/blog/tag/bootstrap-sampling/>), BOOTSTRAPPING (<https://www.analyticsvidhya.com/blog/tag/bootstrap/>), HMISC PACKAGE (<https://www.analyticsvidhya.com/blog/tag/hmisc-package/>), IMPUTE MISSING VALUES (<https://www.analyticsvidhya.com/blog/tag/impute-missing-values/>), IRIS DATA (<https://www.analyticsvidhya.com/blog/tag/iris-data/>), MI PACKAGE



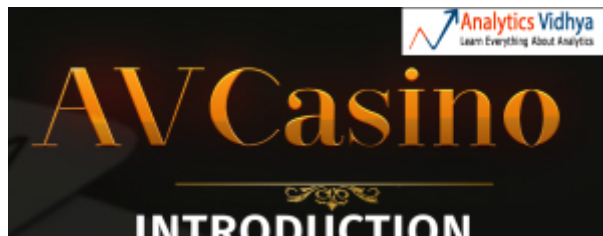
(<https://www.analyticsvidhya.com/blog/tag/mice-package/>), MICE PACKAGE (<https://www.analyticsvidhya.com/blog/tag/mice-package/>), MISSFOREST PACKAGE (<https://www.analyticsvidhya.com/blog/tag/missforest-package/>), MULTIPLE IMPUTATION (<https://www.analyticsvidhya.com/blog/tag/multiple-imputation/>), OUT OF BAG ERROR (<https://www.analyticsvidhya.com/blog/tag/out-of-bag-error/>), PREDICTIVE MEAN MATCHING (<https://www.analyticsvidhya.com/blog/tag/predictive-mean-matching/>)



(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

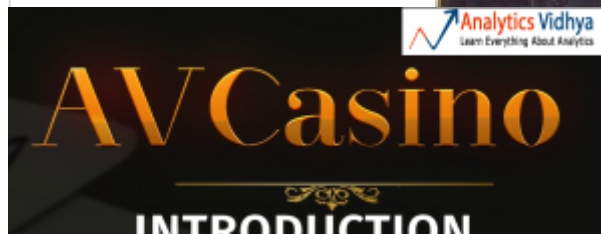
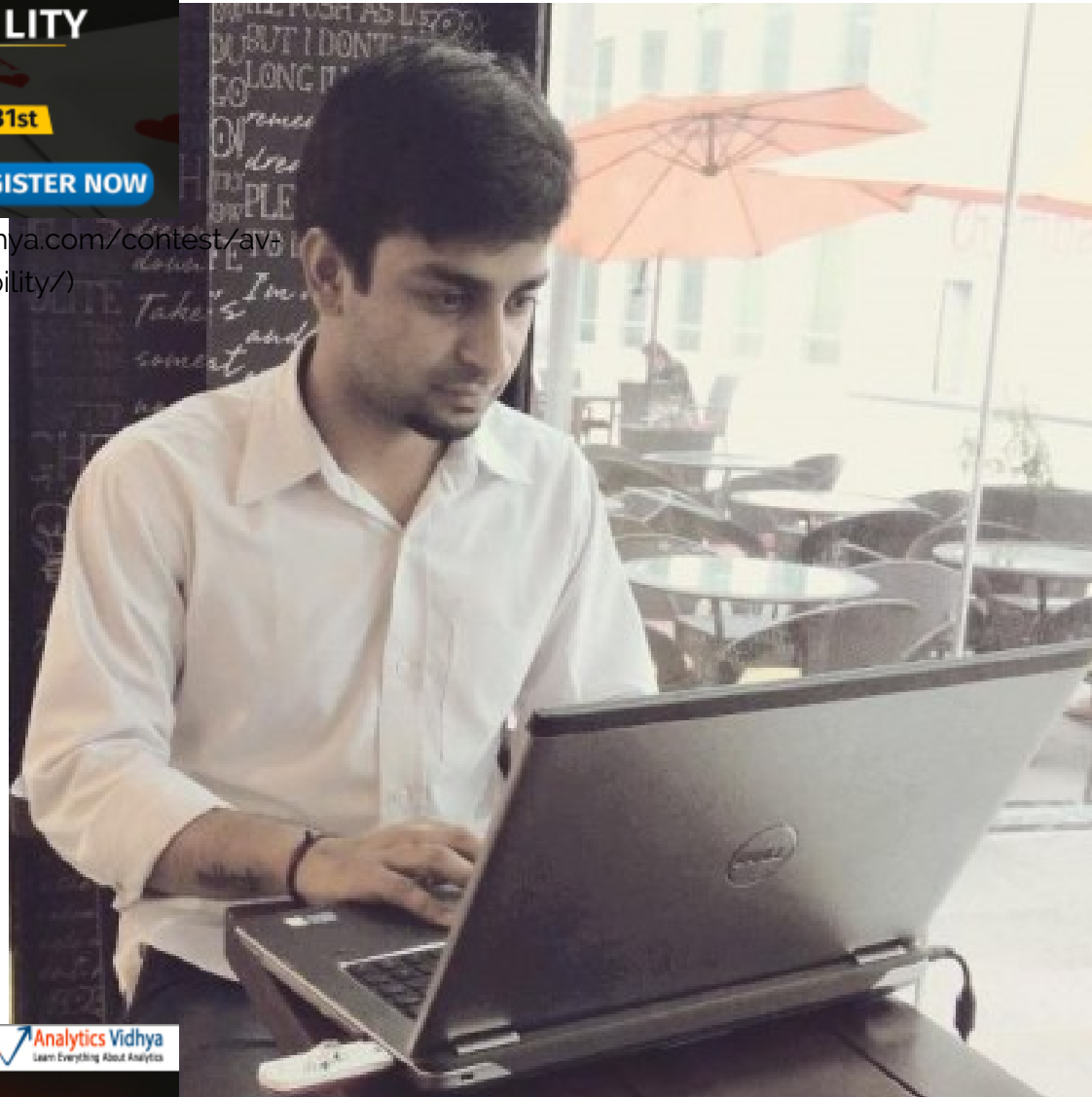
Next Article

10 Questions R Users always ask while using ggplot2 package
(<https://www.analyticsvidhya.com/blog/2016/03/questions-ggplot2-package-r/>)





(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)



(<https://www.analyticsvidhya.com/blog/author/manish-saraswat/>)

Author



(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

<https://www.analyticsvidhya.com/blog/author/manish-saraswat/>)

world. Knowledge is the most powerful asset one can build. It builds up like compound interest. I care about animals, unprivileged people, sharing knowledge, health and books. R, Data Science and Machine Learning keep me busy. Try. Bleed. Succeed.

(https://www.analyticsvidhya.com/Manish_Saraswt) [in](https://in.linkedin.com/in/saraswatmanish) (<https://in.linkedin.com/in/saraswatmanish>)

32 COMMENTS



Surya1987 says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106559#respond>)

MARCH 4, 2016 AT 7:15 AM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-106559>)

Hi Manish, thanks for spending your precious time in writing this nice article. I have one doubt whether transformation has to be done after or before imputing missing values. Secondly is there any method to impute outliers. ^



Manish Saraswat says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106563#respond>)

MARCH 4, 2016 AT 8:26 AM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-106563>)

Hi Surya

In case of Amelia, if the data does not have multivariate normal distribution, transformation is required. Alternatively, you can use [Amelia](#) package. It also uses predictive mean matching, bootstrapping and addition regression methods.





REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=106569#RESPOND)

(https://data.analyticsvidhya.com/forums/106569/introduction-to-probability/) (https://data.analyticsvidhya.com/forums/106569/introduction-to-probability/)

casino-introduction-to-probability/)

Thanks Manish for an excellent article. . For a feature, how much % of values if missing should be considered for imputation ? What I mean is – if a feature has values in 5-10 % of total rows – it is good to drop the feature. Please correct my understanding if I am wrong.

Thanks again!

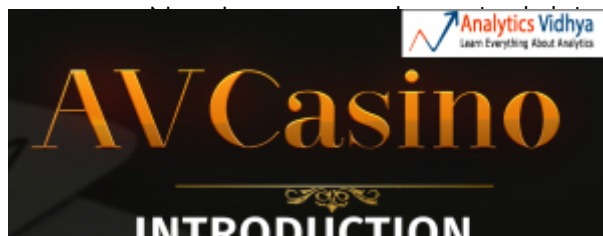


Surya Prakash says:

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=106589#RESPOND)

MARCH 4, 2016 AT 6:14 PM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-106589)

```
newdata<-read.csv(file="C:\\Users\\e885735\\Desktop\\Prakash\\train_u6lujuX.csv",head=TRUE,sep=";",stringsAsFactors =
TRUE,na.strings=c("", "NA", "-", "?"))
newdata1<-na.omit(newdata)
newdata$Credit_History<-as.factor(newdata$Credit_History)
install.packages("missForest")
library(missForest)
newdata.imp<-missForest(newdata[c(2,3,4,5,6,7,8,9,10,11,12,13)])
```



accuracy. However I got the below error

```
imp$ximp,newdata,newdata1)
```



non-numeric argument to binary operator

```

3: In as.character(as.matrix(ximp[, t.ind])) != as.character(as.matrix(xtruel, :
longer object length is not a multiple of shorter object length
4: In as.character(as.matrix(ximp[, t.ind])) != as.character(as.matrix(xtruel, :
longer object length is not a multiple of shorter object length
5: In as.character(as.matrix(ximp[, t.ind])) != as.character(as.matrix(xtruel, :
longer object length is not a multiple of shorter object length

```



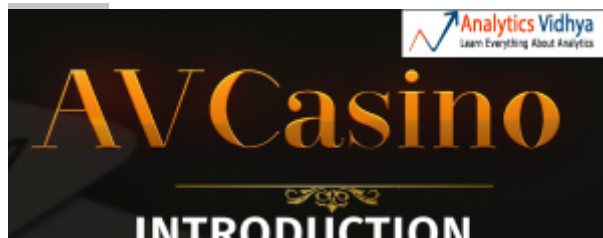
Manish Saraswat says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106738#respond>)

MARCH 7, 2016 AT 12:59 AM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-106738>)

Hi Surya

The error "Longer object length is not a multiple of shorter object length" pops up when one tries to compare two data frames / vectors / arrays of unequal dimensions or sizes. In your case, newdata1 has only 641 observations as compared to newdata which has 981 observations. Since we don't have complete data, it would be difficult to check the accuracy of imputed values. Alternatively, OOB error is also a good estimate of error accuracy. You can always check OOB error using `newdata.imp$OOBerror`



REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106876#respond>)

[ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-106876](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-106876))

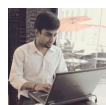


Nalin Pasricha says:
(https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

It should be the arguments in mixError function. In the example which you have provided you have however in my case newdata contains missing values. newdata.imp\$xim is the imputed dataset. and argument in mixError function.

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=106685#RESPOND)
MARCH 6, 2016 AT 5:01 AM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-106685)

great article Manish. I've been using some of these packages for a while but I wasn't aware of many of the nuances you pointed out. Really useful.



Manish Saraswat says:

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=106739#RESPOND)
MARCH 7, 2016 AT 1:00 AM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-106739)

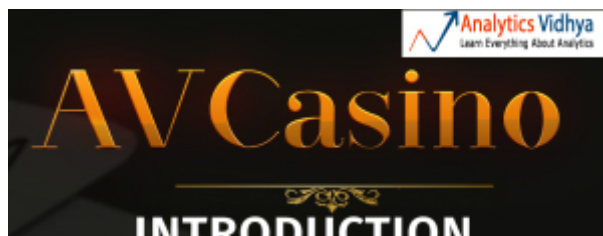
Thanks Nalin.



Maruthi says:

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=106783#RESPOND)
MARCH 7, 2016 AT 7:09 AM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-106783)

Very good information Manish. Could you please throw light on similar methods along with outlier detection in python also?



REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=106928#RESPOND)
ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-106928)



ul, but, I've a problem, this command isn't ok

have class 'mira'

Manish Saraswat says: (https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/) MARCH 10, 2016 AT 6:56 AM (https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-106980) REPLY (https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106980#respond)

Hi Luiz

Generally, this error doesn't pops up. But you can solve it like this:

```
>combine <- pool(as.mira(fit))
```



Vishwa says:

REPLY (https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=109415#respond)

APRIL 13, 2016 AT 6:33 PM (https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-109415)

Hi, I tried combine<-pool(as.mira(fit)) and got this message: Error in pool(as.mira(fit)) : Object has no coef() method.



geeta chhabra says:

REPLY (https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=110137#respond)

APRIL 28, 2016 AT 8:40 AM (https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-110137)

Hi Manish

After using combine<-pool(as.mira(fit))



as no coef() method



REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=110173#respond>)
[ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-110173](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-110173))

as no coef() method.
(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)
Please sort this out
Thanks



Neeraj Agrawal says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=114126#respond>)

JULY 27, 2016 AT 11:19 AM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-114126>)

Hi Manish,

I got the same error. But instead of iris.mis, I used data = imputed_data. If the input of with() is not mids object, it is invoking base with() function.

Please clarify if I am doing anything wrong.

Thanks,



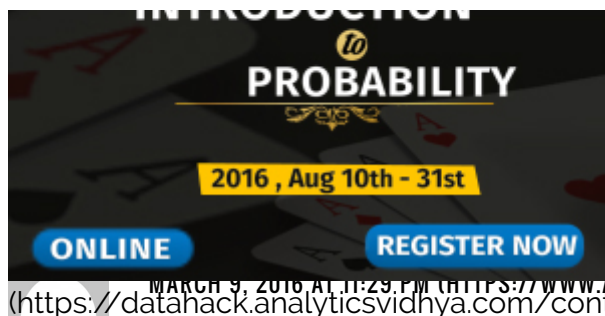
Manimaran says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=114429#respond>)

[ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-114429](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-114429))



pls help me out of this.



REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/?REPLYTOCOM=106945#RESPOND](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106945#respond))
 (https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

I tried to impute with

```
df2dosimputados<-aregImpute(~.,data= df2dosPrestamoslimpio,n.impute=5)
```

my aim is to impute all my vars, but I obtain this error

Error in terms.formula(formula, specials = "I") :
 '.' in formula and no 'data' argument

Do you have any idea to impute all my data frame?

Thanks

.



Manish Saraswat says:

REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/?REPLYTOCOM=106979#RESPOND](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=106979#respond))
[ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-106979](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-106979))





azul77 says:

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)
 MARCH 30, 2016 AT 7:03 PM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-108683>)

Hi Manish thanks a lot, You're right,

I separated my dataframe in two, the first one with columns with nulls values and the second with not nulls values in the columns.

I applied the method to the columns with NA's, but now I have a new trouble, when I check the results, for example `dataframe$imputed$Ultimosmovimientos[,1]`, I only can see the imputed values but not all mi columns values.

Maybe that's not a problem with only one column, I think I could merge the values manually, but I have about 50 columns, so my question is, Do you have any advice to "merge" the imputed values with the values that weren't being imputed.

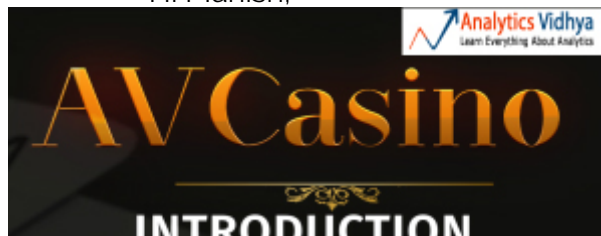
Thanks



Pallavi says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=107487#respond>)
 MARCH 16, 2016 AT 10:58 AM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-107487>)

Hi Manish,



imputing missing values for various projects. And I always used imputation based on some logic.

at we can measure the error in imputation, It made me think how can we check the error.

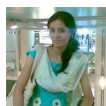
has missing values and we are trying to fill up the data using appropriate logic to predict what's the
 ever know if the prediction is correct. But since we are measuring the accuracy of imputation, I am
 the accuracy against?



REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=107572#RESPOND)
ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-107572)

(https://datahack.analyticsvidhya.com/contest/casino-introduction-to-probability/) You are absolutely right. Missing values don't allow us to check their accuracy (predicted). However, missForest provides us out of bag error estimate. Stekhoven and Bühlmann [2011] showed that this estimate produces an appropriate representation of the true imputation error. Least is desirable.

Alternatively, you can use a long method too. Make different models by using multiple techniques (missForest, Hmisc, mean, median)for missing values imputation. I did it one day. I made 4 different models and found Hmisc performed better & faster.



priyatamil (<http://www.thinkittraining.in/salesforce>) says:

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=107562#RESPOND)
MARCH 17, 2016 AT 5:30 AM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-107562)

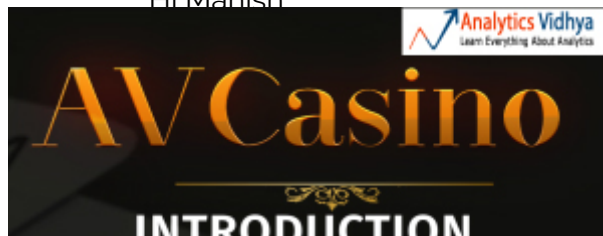
You are said another one valuable information,about the reports was really very great.After refer that post i get new more information,thanks for your valuable support to share that post.



Sudhakar T says:

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/?REPLYTOCOM=108311#RESPOND)
MARCH 26, 2016 AT 6:23 AM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPUTING-MISSING-VALUES/#COMMENT-108311)

Hi Manish



as you describe above. It's new for me. In my case, i am facing a issue related imputation in my data for variables and observation near 15000. In data set, half of predictor variables show completed second half predictor variables show 97% missing cases. Can you recommend which method is



(<https://data.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

can you please help me with getting "iris" data set used in above example...

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=108315#respond>)
[ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-108315](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-108315))



Doug Dame says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=108363#respond>)

MARCH 27, 2016 AT 12:34 AM (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-108363>)

Very interesting article, much thanks.

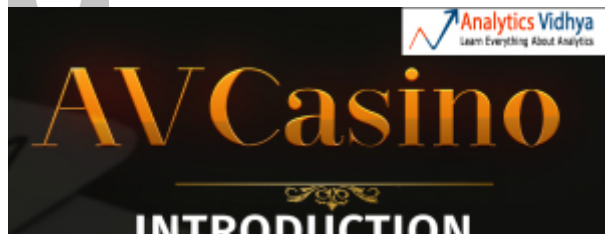
In this case, since you created the missing values in the IRIS dataset yourself, "ground truth" is available. And thus you could show exactly how accurate each of the various methods' imputations were. ^

Doesn't mean those same results would necessarily extrapolate to other datasets, especially ones with more complicated data, but it'd be fun to see !



Rahul says:

REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=109888#respond>)



[ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-109888](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-109888))

it keeps running out of memory.

to ram.



size 34.8 Gb
:
Mb: see help(memory.size)
:
(https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)
3: In rep.int(c(1, numeric(n)), n - 1L) :
Reached total allocation of 8072Mb: see help(memory.size)
4: In rep.int(c(1, numeric(n)), n - 1L) :
Reached total allocation of 8072Mb: see help(memory.size)

The data has about 70K obs. of 12 variables. What should I do?

Thanks



Vamshi Krishna says:

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/?REPLYTOCOM=110359#RESPOND)

MAY 3, 2016 AT 11:38 AM (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-110359)

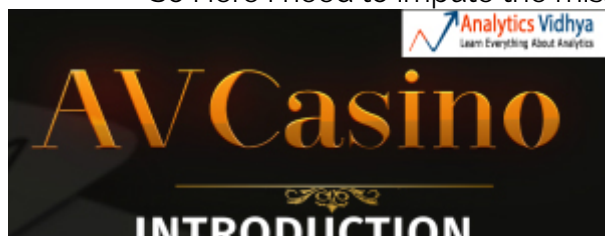
Hi all,

I'm Working on a retail project , I need missing value imputation code in R.

The Dataset is like.

Manufacture > Sub Category > Brand > Sub Brand> Units..

So Here I need to impute the missing values by Manufacture > Sub Category > Brand > Sub Brand wise.





REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/?REPLYTOCOM=111159#RESPOND](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=111159#respond))
[ALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-111159](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-111159))

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)



Avinash says:

REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/?REPLYTOCOM=111690#RESPOND](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=111690#respond))

JUNE 1, 2016 AT 12:33 PM ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-111690](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-111690))

Hello manish

Like in using missForest model using data set of Big Mart Sale, I separated the numerical variables and applied missForest after which when I am trying to use cbind to join the numerical and factor variables to form the original data set it is showing "Error in as.data.frame.default(x[[i]], optional = TRUE, stringsAsFactors = stringsAsFactors) : cannot coerce class ""missForest"" to a data.frame"

I even tried as.data.frame() to change class but it didn't worked out



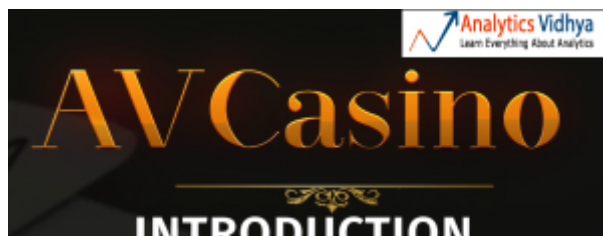
Mudit says:

REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/?REPLYTOCOM=112149#RESPOND](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=112149#respond))

JUNE 12, 2016 AT 1:04 PM ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-112149](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-112149))

Hi,

After running the code using MICE package for imputation this is the error i get



ed_Data1,2)
 (table) :
 d for function 'complete' for signature ""mids""



REPLY (<https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/?replytocom=113823#respond>)
[ALYTICSVIDHYA.COM/BLOG/2016/03/TUTORIAL-POWERFUL-PACKAGES-IMPETING-MISSING-VALUES/#COMMENT-113823](https://www.analyticsvidhya.com/blog/2016/03/tutorial-powerful-packages-imputing-missing-values/#comment-113823))

([https://data.analyticsvidhya.com/contest/av-](https://data.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

casino-introduction-to-probability/

(<http://www.mediafire.com/download/i2nc2di5p4nfbsl/hmisc2.csv>)

It has all of data types.

I use Hmisc package to handle missing values.

My code is:

```
iris.mis=read.csv2("G:\\Thanh Phuong xlsl\\hmisc2.csv", sep=";", na.strings = "na", header=TRUE)
```

```
library(Hmisc)
```

```
impute_arg <- aregImpute(~ weight + oral + gcs + oi + ivdu + csw + previousTB + pulmonaryTB + TBMgrade + disability.base+  
disability.2mo+ cd4count+ cd4.2mo+ hivrna.base+ hivrna.2mo, data = iris.mis, n.impute = 5)
```

and i have a notice:

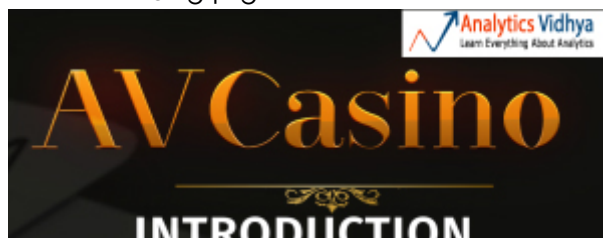
Iteration 1

fewer than 3 unique knots. Frequency table of variable:

x

1 2 3

61 54 15



```
ms, nk = nk, inclx = TRUE) :
```

```
, nk = nk, inclx = TRUE) :
```

with default algorithm.



h 3 knots
, nk = nk, inclx = TRUE) :
values of x. knots set to 1 interior values.

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

LEAVE A REPLY

Connect with:



(<https://www.analyticsvidhya.com/wp-login.php?>

action=wordpress_social_authenticate&mode=login&provider=Facebook&redirect_to=https%3A%2F%2Fwww.analyticsvidhya.com%2Fblog%2F2016-03-24/tutorial-on-5-powerful-packages-imputing-missing-values%2F)

Your email address will not be published.

Comment










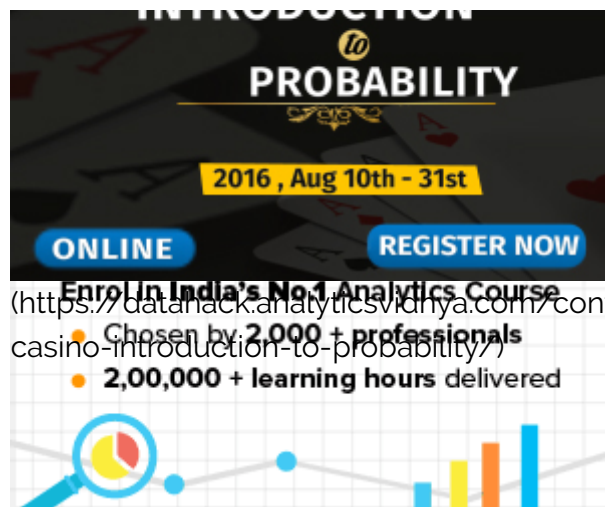
(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

SUBMIT COMMENT

TOP AV USERS

Rank	Name	Points
1	 SRK (https://datahack.analyticsvidhya.com/user/profile/SRK)	5378
2	 Aayushmnit (https://datahack.analyticsvidhya.com/user/profile/aayushmnit)	4818
3	 Nalin Pasricha (https://datahack.analyticsvidhya.com/user/profile/Nalin)	4407
	 datahack.analyticsvidhya.com/user/profile/Rohan Rao	4353
	 datahack.analyticsvidhya.com/user/profile/binga	3371

More Rankings (<http://datahack.analyticsvidhya.com/users>)



(<http://www.greatlearning.in/great-lakes-pgpba?>

([https://datahack.analyticsvidhya.com/contest/av-](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

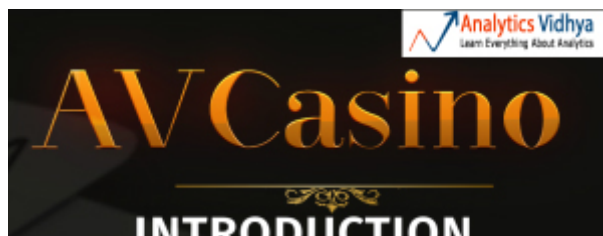
[casino-introduction-to-probability/](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/))

• 2,00,000+ learning hours delivered

utm_source=avm&utm_medium=avmbanner&utm_campaign=pgpba)

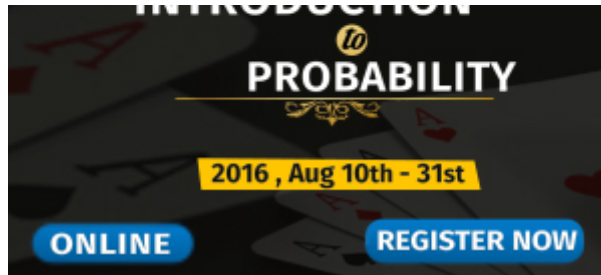
POPULAR POSTS

- A Complete Tutorial to Learn Data Science with Python from Scratch (<https://www.analyticsvidhya.com/blog/2016/01/complete-tutorial-learn-data-science-python-scratch-2/>)
- 7 Types of Regression Techniques you should know! (<https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/>)
- Essentials of Machine Learning Algorithms (with Python and R Codes) (<https://www.analyticsvidhya.com/blog/2015/08/common-machine-learning-algorithms/>)
- A Complete Tutorial on Tree Based Modeling from Scratch (in R & Python) (<https://www.analyticsvidhya.com/blog/2016/04/complete-tutorial-tree-based-modeling-scratch-in-python/>)



Modeling in R (<https://www.analyticsvidhya.com/blog/2015/12/complete-tutorial-time-series->

Forecast (with Codes in Python) (<https://www.analyticsvidhya.com/blog/2016/02/time-series->

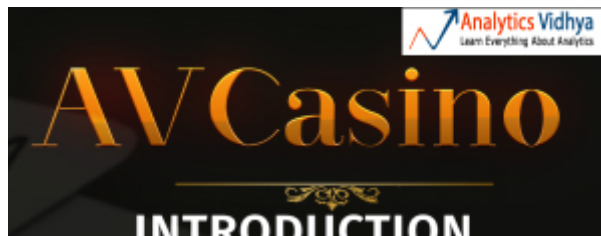


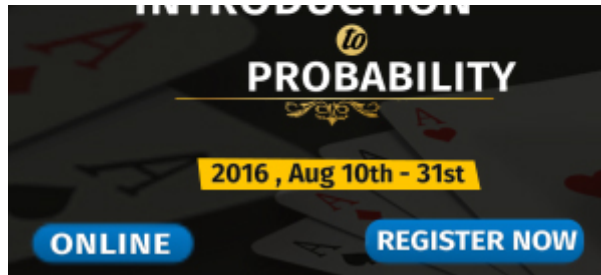
thon (using BeautifulSoup) (<https://www.analyticsvidhya.com/blog/2015/10/beginner-guide-web->

for Data Manipulation (<https://www.analyticsvidhya.com/blog/2016/01/12-pandas-techniques->

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

FEATURED VIDEO

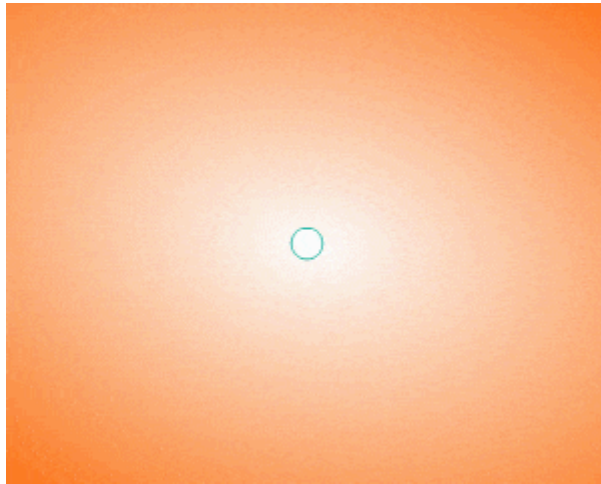




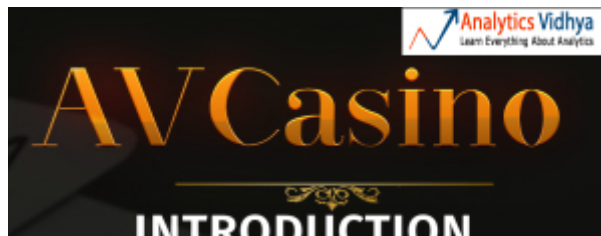
Weeks at Marketing Conclave 2014



(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)



(<http://imarticus.org/sas-online>)





[https://datahack.analyticsvidhya.com/contest/av-](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

[casino-introduction-to-probability/](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/) <https://www.analyticsvidhya.com/blog/2016/08/bringing-analytics-into-indian-film-industry-with-back-tracing-algorithm/>



Analytics into Indian Film Industry with Back Tracing Algorithm (<https://www.analyticsvidhya.com/blog/2016/08/bringing-analytics-into-indian-film-industry-with-back-tracing-algorithm/>)

GUEST BLOG , AUGUST 22, 2016



(<https://www.analyticsvidhya.com/blog/2016/08/industry-insight-fighting-cyber-fraud-with-analytics/>)

(<https://www.analyticsvidhya.com/blog/2016/08/industry-insight-fighting-cyber-fraud-with-analytics/>)

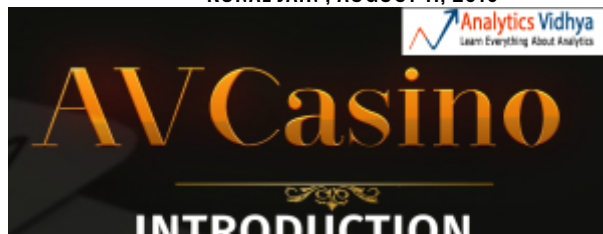
GUEST BLOG , AUGUST 15, 2016



(<https://www.analyticsvidhya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/>)

hya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/)

KUNAL JAIN , AUGUST 11, 2016



Beginners Guide to Topic Modeling in Python
(<https://www.analyticsvidhya.com/blog/2016/08/beginners-guide-to-topic-modeling-in-python/>)

**Br
ng
ng**

Industry Insight - Fighting Cyber Fraud with Analytics



Launch of AV Casino - An Introduction to Probability
(<https://www.analyticsvidhya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/>)



(<http://www.edvancer.in/course/cbap?>

([https://datahack.analyticsvidhya.com/contest/av-](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

[casino-introduction-to-probability/](https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/))

[utm_source=AV&utm_medium=AVads&utm_campaign=AVadsnonfc&utm_content=cbapavad](#))

GET CONNECTED



6,044

FOLLOWERS

(<http://www.twitter.com/analyticsvidhya>)



1,264

FOLLOWERS

(<https://plus.google.com/+Analyticsvidhya>)



18,628

FOLLOWERS

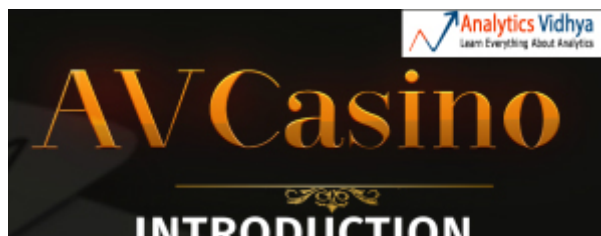
(<http://www.facebook.com/Analyticsvidhya>)

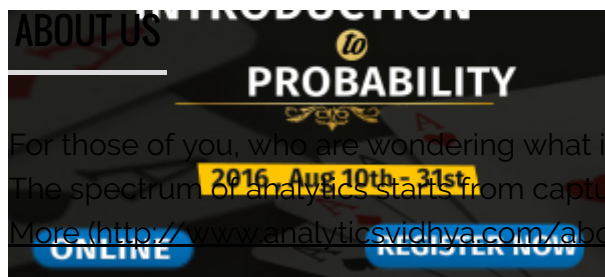


Email

SUBSCRIBE

([http://feedburner.google.com/fb/a/mailverify?](http://feedburner.google.com/fb/a/mailverify?uri=analyticsvidhya)
[uri=analyticsvidhya](http://feedburner.google.com/fb/a/mailverify?uri=analyticsvidhya))





For those of you, who are wondering what is "Analytics Vidhya", "Analytics" can be defined as the science of extracting insights from raw data. The spectrum of analytics starts from capturing data and evolves into using insights / trends from this data to make informed decisions. [Read More \(http://www.analyticsvidhya.com/about-me/\)](http://www.analyticsvidhya.com/about-me/)

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

STAY CONNECTED



6,044

FOLLOWERS

(<http://www.twitter.com/analyticsvidhya>)



1,264

FOLLOWERS

(<https://plus.google.com/+Analyticsvidhya>)



18,628

FOLLOWERS

(<http://www.facebook.com/Analyticsvidhya>)



Email

SUBSCRIBE

(<https://feedburner.google.com/fb/a/mailverify?uri=analyticsvidhya>)

LATEST POSTS



(<https://www.analyticsvidhya.com/blog/2016/08/beginners-guide-to-topic-modeling-in-python/>)

[om/blog/2016/08/beginners-guide-to-topic-modeling-in-python/](https://www.analyticsvidhya.com/blog/2016/08/beginners-guide-to-topic-modeling-in-python/))

SHIVAM BANSAL , AUGUST 24, 2016

Beginners Guide to Topic Modeling in Python
(<https://www.analyticsvidhya.com/blog/2016/08/beginners-guide-to-topic-modeling-in-python/>)



(<https://www.analyticsvidhya.com/blog/2016/08/bringing-analytics-into-indian-film-industry-with-back-tracing-algorithm/>)

**Br
ng
ng**



Analytics into Indian Film Industry with Back Tracing Algorithm (<https://www.analyticsvidhya.com/blog/2016/08/bringing-analytics-into-indian-film-industry-with-back-tracing-algorithm/>)

GUEST BLOG , AUGUST 22, 2016

2016 , Aug 10th - 31st

REGISTER NOW

(<https://www.analyticsvidhya.com/blog/2016/08/industry-insight-fighting-cyber-fraud-with-analytics/>)

[av-casino-introduction-to-probability/](https://www.analyticsvidhya.com/contest/av-casino-introduction-to-probability/)

(<https://www.analyticsvidhya.com/blog/2016/08/industry-insight-fighting-cyber-fraud-with-analytics/>)

GUEST BLOG , AUGUST 15, 2016

Industry Insight - Fighting Cyber Fraud with Analytics



(<https://www.analyticsvidhya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/>)

[hya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/](https://www.analyticsvidhya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/)

KUNAL JAIN , AUGUST 11, 2016

Launch of AV Casino - An Introduction to Probability
(<https://www.analyticsvidhya.com/blog/2016/08/launch-of-av-casino-an-introduction-to-probability/>)

QUICK LINKS

[Home](https://www.analyticsvidhya.com/) (<https://www.analyticsvidhya.com/>)

[About Us](https://www.analyticsvidhya.com/about-me/) (<https://www.analyticsvidhya.com/about-me/>)

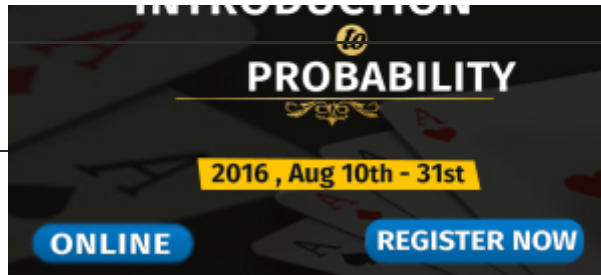
[Our team](https://www.analyticsvidhya.com/about-me/team/) (<https://www.analyticsvidhya.com/about-me/team/>)

[Privacy Policy](https://www.analyticsvidhya.com/privacy-policy/) (<https://www.analyticsvidhya.com/privacy-policy/>)

[Refund Policy](https://www.analyticsvidhya.com/refund-policy/) (<https://www.analyticsvidhya.com/refund-policy/>)

[Terms of Use](https://www.analyticsvidhya.com/terms/) (<https://www.analyticsvidhya.com/terms/>)





© Copyright 2016 Analytics Vidhya

(<https://datahack.analyticsvidhya.com/contest/av-casino-introduction-to-probability/>)

