

## 采用

### 一、判断缺失值

缺失值有两种表示方式：NA和NaN，两者之间的关系是：NaN属于NA，NA不属于NaN。原因是：NaN一般表示数值类型的缺失值，NA表示的数据类型可以多种，比如：数值的缺失值、字符的缺失值等。

#### 1.先用c()函数创建一个向量

```
1 > x<-c(1,NA,NA,2,3)
2 > x
3 [1] 1 NA NA 2 3
```

#### 2.用is.na()和is.nan()函数判断缺失值

```
1 > is.na(x)
2 [1] FALSE TRUE TRUE FALSE FALSE
3 > is.nan(x)
4 [1] FALSE FALSE FALSE FALSE FALSE
```

从结果可以看出：is.na()函数判断出了每个元素是否为缺失值，缺失值返回了真，数值返回了假；is.nan()函数没有判断出缺失值，说明在向量x中不存在NaN这种类型的缺失值，同时也印证了NA不属于NaN。若我们把向量x里的NA改成NaN，结果如下：

```
1 > x<-c(1,NaN,NaN,2,3)
2 > x
3 [1] 1 NaN NaN 2 3
4 > is.na(x)
5 [1] FALSE TRUE TRUE FALSE FALSE
6 > is.nan(x)
7 [1] FALSE TRUE TRUE FALSE FALSE
```

### 二、处理缺失值

```

1 #1.先用c()函数创建一个包含缺失值的向量
2 > x<-c(1,NA,NA,2,3)
3 > x
4 [1] 1 NA NA 2 3
5 #2.取向量中不是缺失值的元素，！的意思是取反
6 > x[!is.na(x)]
7 [1] 1 2 3
8
9 #3.用complete.case()函数取多个向量对应位置都不是缺失值的元素
10 > x<-c(1,NA,NA,2,3)
11 > y<-c(NA,"a","b","c","d")
12 > z<-c(1L,2L,3L,NA,4L)
13 > w<-complete.cases(x,y,z)
14 > w
15 [1] FALSE FALSE FALSE FALSE TRUE
16 #返回结果是逻辑向量，其中只有x和y对应位置都不是缺失值的才会返回TRUE，否则返回FALSE。我们用x[

```

```

1 > x[w]
2 [1] 3
3 > y[w]
4 [1] "d"
5 > z[w]
6 [1] 4

```

## 三、实例

### 1.加载数据集所在的包

```

1 > library(datasets)

```

### 2.用head()函数查看数据集的前6行

```

1 head(airquality)
2   Ozone Solar.R Wind Temp Month Day
3 1    41    190  7.4  67     5   1

```

4	2	36	118	8.0	72	5	2
5	3	12	149	12.6	74	5	3
6	4	18	313	11.5	62	5	4
7	5	NA	NA	14.3	56	5	5
8	6	28	NA	14.9	66	5	6

### 3.用complete.cases()函数选择在变量

```

1 g<-complete.cases(airquality)
2 > g
3 [1] TRUE TRUE TRUE TRUE FALSE FALSE TRUE TRUE TRUE FALSE
4 [11] FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
5 [21] TRUE TRUE TRUE TRUE FALSE FALSE FALSE TRUE TRUE TRUE
6 [31] TRUE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE
7 [41] TRUE FALSE FALSE TRUE FALSE FALSE TRUE TRUE TRUE TRUE
8 [51] TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
9 [61] FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE
10 [71] TRUE FALSE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE
11 [81] TRUE TRUE FALSE FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
12 [91] TRUE TRUE TRUE TRUE TRUE FALSE FALSE FALSE TRUE TRUE
13 [101] TRUE FALSE FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE
14 [111] TRUE TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE FALSE TRUE
15 [121] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
16 [131] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
17 [141] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE FALSE
18 [151] TRUE TRUE TRUE
19

```

### 4.选择该数据集不包含缺失值记录的前10行，且每个变量都要

```

1 > airquality[g,][1:10,]
2   Ozone Solar.R Wind Temp Month Day
3 1    41    190  7.4  67    5    1
4 2    36    118  8.0  72    5    2
5 3    12    149 12.6  74    5    3
6 4    18    313 11.5  62    5    4
7 7    23    299  8.6  65    5    7
8 8    19     99 13.8  59    5    8

```

9	9	8	19	20.1	61	5	9
10	12	16	256	9.7	69	5	12
11	13	11	290	9.2	66	5	13
12	14	14	274	10.9	68	5	14

## 四、识别缺失值的模式

```
1 install.packages("mice")
2 library(mice)
```

```
1 #用md.pattern()函数查看数据缺失值的分布
2 md.pattern(data)
```

## 五、处理缺失值

### 1.行删除法

用na.omit()函数删除不完整观测

```
1 library(datasets)
2 > data<-head(airquality)
3 > data
4   Ozone Solar.R Wind Temp Month Day
5 1    41    190  7.4  67     5   1
6 2    36    118  8.0  72     5   2
7 3    12    149 12.6  74     5   3
8 4    18    313 11.5  62     5   4
9 5    NA     NA 14.3  56     5   5
10 6    28     NA 14.9  66     5   6
11 > newdata<-na.omit(data) #用na.omit()函数删除不完整观测
12 > newdata
13   Ozone Solar.R Wind Temp Month Day
14 1    41    190  7.4  67     5   1
15 2    36    118  8.0  72     5   2
16 3    12    149 12.6  74     5   3
17 4    18    313 11.5  62     5   4
```

## 2.多重插补法

多重插补法（MI）是一种基于重复模拟的处理缺失值的方法，多重插补是从一个包含缺失值的数据集中生成一组完整的数据集。每个模拟数据集中，缺失数据将使用蒙特卡洛方法来填补。

```
1 library(datasets)
2 data<-head(airquality)
3 data
4
5 data1<-mice(data,m=6)
6 #用mice()函数从包含缺失数据的数据框开始，进行6重插补，即生成6个完整数据集
7 #mice()函数的第一个参数data为数据，第二个参数m为要返回的完整数据集的个数
8
9 fit<-with(data1,lm(Solar.R~Wind+Temp+Month+Day+Ozone))
10 fit
11 #用with()函数依次对6个完整数据集分别应用lm()模型
12 #结果分别返回6个完整数据集的回归结果
13
14 jiegua<-pool(fit)
15 jiegua
16 summary(jiegua)
17 #用pool()函数汇总回归结果
18
19 result<-complete(data1,action=3)
20 result
21 #选择第三个插补数据集作为结果
```