

## Missing value handling in R

# MISSING VALUE HANDLING IN R



# Missing value handling in R

## In this session

1. Replace missing values with the mean
2. Replace missing values with the median
3. Replace missing values with an interpolated estimate
4. Replace missing values with a constant
5. Replace missing values using imputation
6. Replace missing values with a missing rank
7. Replace missing values with a dummy
8. Replace missing values with 0
9. Create an indicator variable for "missing."
10. Replace missing values with a string
11. Add an indicator variable showing which strings are considered "missing."
12. Delete columns that are missing too many values to be useful
13. Delete rows that are missing critical values

# Missing value handling in R




























We need data that is:

- Relevant
- Connected
- Accurate
- Enough to work with



# Missing value handling in R

## Example of missing values dataset

Column 0	age	years_seniority	income	parking_space	attending_party	entree	pets	emergency_contact
								
Tony	48	27		1	5	shrimp		Pepper
Donald	67	25	86	10	2	beef		Jane
Henry	69	21	95	6	1	chicken	62	Janet
Janet	62	21	110	3	1	beef		Henry
Nick		17		4				
Bruce	37	14	63		1	veggie		NA
Steve	83		77	7	1	chicken		n/a
Clint	27	9	118	9		shrimp	3	None
Wanda	19	7	52	2	2	shrimp		empty
Natasha	26	4	162	5	3			-
Carol		3	127	11	1	veggie	1	""
Mandy	44	2	68	8	1	chicken		null

# Missing value handling in R



Too many missing data == Swiss cheese

# Missing value handling in R

Example of missing values dataset CSV file

missing\_values.csv

	A	B	C	D	E	F	G	H	I
1		age	years_seniority	income	parking_space	attending_party	entree	pets	emergency_contact
2	Tony	48	27		1	5	shrimp		Pepper
3	Donald	67	25	86	10	2	beef		Jane
4	Henry	69	21	95	6	1	chicken	62	Janet
5	Janet	62	21	110	3	1	beef		Henry
6	Nick		17		4				
7	Bruce	37	14	63		1	veggie		NA
8	Steve	83		77	7	1	chicken		n/a
9	Clint	27	9	118	9		shrimp	3	None
10	Wanda	19	7	52	2	2	shrimp		empty
11	Natasha	26	4	162	5	3			_
12	Carol		3	127	11	1	veggie	1	""
13	Mandy	44	2	68	8	1	chicken		null

# Missing value handling in R

## R Studio

The screenshot shows the RStudio interface with a data table, console output, and a sidebar with documentation.

**Data Table:**

	name	age	years_seniority	income	parking_space	attending_party	entree	pets	emergency_contact
1	Tony	48	27	NA	1	5	shrimp	NA	Pepper
2	Donald	67	25	86	10	2	beef	NA	Jane
3	Henry	69	21	95	6	1	chicken	62	Janet
4	Janet	62	21	110	3	1	beef	NA	Henry
5	Nick	NA	17	NA	4	NA	NA	NA	NA
6	Bruce	37	14	63	NA	1	veggie	NA	NA
7	Steve	83	NA	77	7	1	chicken	NA	n/a
8	Clint	27	9	118	9	NA	shrimp	3	None
9	Wanda	19	7	52	2	2	shrimp	NA	empty
10	Natasha	26	4	162	5	3	NA	NA	-
11	Carol	NA	2	127	11	1	veggie	1	**

Showing 1 to 11 of 12 entries

**Console:**

```

$ years_seniority : int 27 25 21 21 17 14 NA 9 7 4 ...
$ income          : int  NA 86 95 110 NA 63 77 118 52 162 ...
$ parking_space   : int  1 10 6 3 4 NA 7 9 2 5 ...
$ attending_party  : int  5 2 1 1 NA 1 1 NA 2 3 ...
$ entree          : Factor w/ 4 levels "beef","chicken",...: 3 1 2 1 NA 4 2 3 3 NA
...
$ pets           : int  NA NA 62 NA NA NA NA 3 NA NA ...
$ emergency_contact: Factor w/ 11 levels "\"\"","_","empty",...: 11 5 6 4 NA 8 7 9 3 2 ...
> # Check the data for missing values
> sapply(dat, function(x) sum(is.na(x)))
      name      age years_seniority      income      parking_
space      0      2              1              2
1
  attending_party      entree      pets emergency_contact
2              2              9              1
> View(dat)
>
  
```

**Environment:** dat 12 obs. of 9 ...

**Files Plots Packages Help Viewer**

**R: Arithmetic Mean** Find in Topic

**mean (base)** R Documentation

**Arithmetic Mean**

**Description**

Generic function for the (trimmed) arithmetic mean.

**Usage**

```
mean(x, ...)
```

## Default S3 method:  
mean(x, trim = 0, na.rm = F.

**Arguments**

x An R object. Currently there are methods for numeric/logical vectors and [date](#), [date-time](#) and [time](#)

# Missing value handling in R

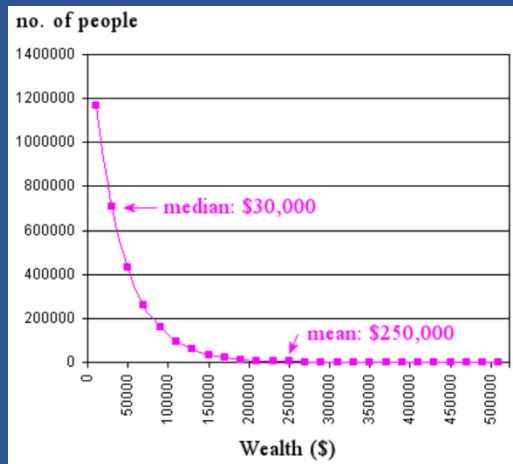
## General commands

```
1 rm(list = ls()) # clear work space
2 setwd("d:/temp") # set current work directory
3 sessionInfo() # get session information
4 installed.packages() # list installed packages
5 # import data file
6 dat <- read.csv("missing_values.csv", na.strings = "")
7 str(dat) # show data frame structure
8 # Check the data for missing values
9 sapply(dat, function(x) sum(is.na(x)))
```



# Missing value handling in R

Replace missing values with the mean



Sample Mean	Population Mean
$\bar{x} = \frac{\sum x}{n}$	$\mu = \frac{\sum x}{N}$

where  $\sum x$  is sum of all data values

$N$  is number of data items in population

$n$  is number of data items in sample

```

1  # Replace missing values with the mean
2  # column = age
3  # Missing values type = distributed
4  # Formal name = Missing Completely at Random (MCAR)
5  rm(list = ls()) # clear work space
6  setwd("d:/temp") # set current work directory
7  # import data file
8  dat <- read.csv("missing_values.csv", na.strings = "")
9  dat$age.mean <- ifelse(is.na(dat$age),
10 |                         mean(dat$age, na.rm = TRUE),
11 |                         dat$age)

```

# Missing value handling in R

Replace missing values with the median

$$\text{Median} = l + \frac{h}{f} \left( \frac{N}{2} - c \right)$$

Where:

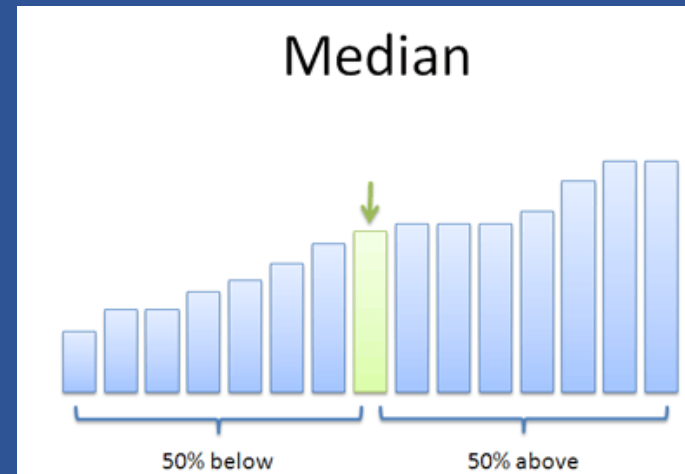
$l$  = lower class boundary of the median class

$h$  = Size of the median class interval

$f$  = Frequency corresponding to the median class

$N$  = Total number of observations i.e. sum of the frequencies

$c$  = Cumulative frequency preceding median class.

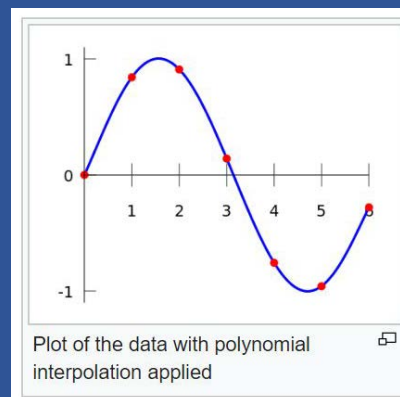


```
1 # Replace missing values with the median
2 # column = age
3 # • Another justifiable way to handle missing-at-random data
4 rm(list = ls()) # clear work space
5 setwd("d:/temp") # set current work directory
6 # import data file
7 dat <- read.csv("missing_values.csv", na.strings = "")
8 dat$age.mean <- ifelse(is.na(dat$age),
9                        median(dat$age, na.rm = TRUE), dat$age)
```

# Missing value handling in R

Replace missing values with an interpolated estimate

```
1 # Replace missing values with an interpolated estimate
2 # column = years_seniority
3 # the values in this column, years seniority, is ordered,
4 # greatest to least. This structure can be exploited by
5 # interpolating the missing value. This approach is very effective
6 # when it is appropriate, usually with time - series data.
7 rm(list = ls()) # clear work space
8 setwd("d:/temp") # set current work directory
9 # import data file
10 dat <- read.csv("missing_values.csv", na.strings = "")
11 dat$senior <- ifelse(is.na(dat$years_seniority), 11.5,
12                     dat$years_seniority)
```



# Missing value handling in R

Replace missing values with a constant

```
1 # Replace missing values with a constant
2 # column = income
3 # Missing values are Missing Not at Random(MNAR)
4 # Those with very high incomes preferred not to state them
5 # Make a reasonable guess for what "high" means and fill
6 # in the blanks. It will still be inaccurate, but better than blank
7
8 rm(list = ls()) # clear work space
9 setwd("d:/temp") # set current work directory
10 # import data file
11 dat <- read.csv("missing_values.csv", na.strings = "")
12 dat$income4 <- ifelse(is.na(dat$income), 250, dat$income)
```



# Missing value handling in R

Replace missing values using imputation MICE

```
1  # Column = years_seniority, age, income
2  # Method = multivariate imputation by chained equations(MICE)
3  # MICE is the method of choice for complex incomplete data
4  # Good for missing data in more than one variable
5  # powerful when features are somewhat related
6
7  rm(list = ls())
8  dat <- read.csv("missing_values.csv", na.strings = "")
9  # A fast, consistent tool for working with data frame
10 install.packages("dplyr")
11 library(dplyr) # attach packages
12 # Make new dataframe with column years_seniority,
13 # age, income with column type = numeric
14 dat <- dat %>% # make new 3 columns with type numeric
15   mutate(
16     senior1 = as.numeric(years_seniority),
17     age1 = as.numeric(age),
18     income1 = as.numeric(income)
19   )
```

# Missing value handling in R

Replace missing values using imputation MICE

```
21 # Replace missing values using imputation MICE
22 keep <- c("senior1", "age1", "income1")
23 #drop all columns but keep 3 col
24 dat <- dat[, keep, drop = FALSE]
25 install.packages('mice')           # standard command
26 library(mice)                      # standard command
27 init = mice(dat, maxit = 0)         # standard command
28 meth = init$method                  # standard command
29 predM = init$predictorMatrix
30 # Bayesian linear regression (การวิเคราะห์ถดถอยเชิงเส้นแบบเบย์ส์)
31 meth[c("senior1")] = "norm"
32 Predictive mean matching
33 meth[c("age1")] = "pmm" #
34 meth[c("income1")] = "pmm"
35 # Replace missing values using imputation MICE
36 set.seed(103) # seed for pseudo random number generator
37 imputed = mice(dat, method = meth, predictorMatrix = predM, m = 5)
38 imputed <- complete(imputed)
```

# Missing value handling in R

Replace missing values with a missing rank

```
1 # Replace missing values with a missing rank
2 # Column = parking_space
3 # Our knowledge of how parking spaces are numbered,
4 # let us make a guess here
5 # All the space numbers from 1 - 11
6 # Missing one might be 12
7 rm(list = ls())
8 dat <- read.csv("missing_values.csv", na.strings = "")
9 dat$park <- ifelse(is.na(dat$parking_space), 12,
10                  dat$parking_space)
```

# Missing value handling in R

Replace missing values with a dummy

```
1  # eplace missing values with a dummy
2  # Column = parking_space
3  # Filling in a dummy value
4  # Clearly different from actual values
5  # Such as a negative rank
6  # Used to indicate that the feature is not applicable
7
8  rm(list = ls())
9  dat <- read.csv("missing_values.csv", na.strings = "")
10 dat$park <- ifelse(is.na(dat$parking_space), -99,
11                   dat$parking_space)
```



# Missing value handling in R

Replace missing values with 0

```
1 # eplace missing values with 0
2 # Column = attending_party
3 # A missing numerical value can mean zero.
4 # In the case of an RSVP, invitees who are not planning
5 # to attend sometimes neglect to respond, but guests
6 # planning to attend are more likely to.
7 # In this case, filling in missing blanks with a zero is reasonable
8 rm(list = ls())
9 dat <- read.csv("missing_values.csv", na.strings = "")
10 dat$party <- ifelse(is.na(dat$attending_party), 0, dat$attending_party)
```

# Missing value handling in R

Create an indicator variable for "missing"

```
1  # Create an indicator variable for "missing"
2  # Column = pets
3  # Replacing missing values requires making assumptions.
4  # Whenever your confidence in those assumptions is low,
5  # it is safer to also create a true / false feature
6  # indicating that the value was missing.
7  # This allows many algorithms to learn to weight those differently.
8  rm(list = ls())
9  dat <- read.csv("missing_values.csv", na.strings = "")
10 dat$pet1 <- ifelse(is.na(dat$pet), 0, dat$pets)
11 dat$pet2 <- complete.cases(dat$pets)
```

# Missing value handling in R

## Replace missing values with a string

```
1 # Replace missing values with a string
2 # Column = emergency_contact
3 # Replace NA with "no"
4
5 rm(list = ls())
6 dat <- read.csv("missing_values.csv", na.strings = "")
7 dat$emer <- ifelse(is.na(dat$emergency_contact), 'no',
8                   dat$emergency_contact)
```

## Missing value handling in R

Add an indicator variable showing which strings are considered "missing."

```
1 # Add an indicator variable showing which strings
2 # are considered "missing."
3 # Column = emergency_contact
4 # There are lots of ways to communicate the concept
5 # of "missing" in a string
6 # Replace no, NA, n / a, None, _, "", empty, null with 0
7 # Otherwise = 1
8 rm(list = ls())
9 dat <- read.csv("missing_values.csv", na.strings = "")
10 dat$emer <- ifelse(dat[9] == 'NA' |
11   is.na(dat[9]) | dat[9] == 'n/a' |
12   dat[9] == 'None' | dat[9] == 'empty' |
13   dat[9] == '_' | dat[9] == '' | dat[9] == 'null', 0, 1)
```

# Missing value handling in R

Delete columns that are missing too many values to be useful

```
1 # Column = pets
2 # Is a feature that missing too many values
3 # Not enough information available to make reasonable
4 # assumptions about how to replace the missing values
5 # Best policy = delete the column entirely.
6
7 rm(list = ls())
8 dat <- read.csv("missing_values.csv", na.strings = "")
9 dat <- dat[, !(names(dat) %in% 'pets')]
```

# Missing value handling in R

Delete rows that are missing critical values

```
1 # Delete rows that are missing critical values
2 # Rows that are missing important features can be deleted.
3 # This is particularly useful when you have the luxury of
4 # hand - picking high - quality data
5 # such as when training a model
6 rm(list = ls())
7 dat <- read.csv("missing_values.csv", na.strings = "")
8 dat <- dat[complete.cases(dat),]
```

# Missing value handling in R

## More information on Missing value handling in R

Bayesian linear regression analysis without tears (R)

<https://www.r-bloggers.com/bayesian-linear-regression-analysis-without-tears-r/>

Source code

<https://github.com/laploy/ML/blob/master/Missing%20R%20Script.zip>