

WHAT TO DO WHEN SOME VALUES ARE MISSING

Statistics 407
ISU

1

OUTLINE

- Terminology
- Issues of missingness for multivariate data
- Plotting missings, and describing the distributions of missing vs not missing
- Imputation methods

2

BACKGROUND TERMS

- MCAR: probability that a value is missing does _____ on any other observed or unobserved value.
- MAR: probability that a value is missing _____ only on the _____ variables.
- MNAR: the reason for missing values depends on some _____ information - very difficult analysis.

3

EXAMPLE

Case	X_1	X_2	X_3	X_4	X_5
1	NA	20	1.8	6.4	-0.8
2	0.3	NA	1.6	5.3	-0.5
3	0.2	23	1.4	6.0	NA
4	0.5	21	1.5	NA	-0.3
5	0.1	21	NA	6.4	-0.5
6	0.4	22	1.6	5.6	-0.8
7	0.3	19	1.3	5.9	-0.4
8	0.5	20	1.5	6.1	-0.3
9	0.3	22	1.6	6.3	-0.5
10	0.4	21	1.4	5.9	-0.2

Missing:

_____ of the numbers

_____ of variables

_____ of samples

4

SUMMARY STATISTICS

Case	X_1	X_2	X_3	X_4	X_5
1	NA	20	1.8	6.4	-0.8
2	0.3	NA	1.6	5.3	-0.5
3	0.2	23	1.4	6.0	NA
4	0.5	21	1.5	NA	-0.3
5	0.1	21	NA	6.4	-0.5
6	0.4	22	1.6	5.6	-0.8
7	0.3	19	1.3	5.9	-0.4
8	0.5	20	1.5	6.1	-0.3
9	0.3	22	1.6	6.3	-0.5
10	0.4	21	1.4	5.9	-0.2

Means can be calculated

_____.

Correlations can be
calculated _____.

5

SHADOW MATRIX

Case	X_1	X_2	X_3	X_4	X_5
1	NA	20	1.8	6.4	-0.8
2	0.3	NA	1.6	5.3	-0.5
3	0.2	23	1.4	6.0	NA
4	0.5	21	1.5	NA	-0.3
5	0.1	21	NA	6.4	-0.5
6	0.4	22	1.6	5.6	-0.8
7	0.3	19	1.3	5.9	-0.4
8	0.5	20	1.5	6.1	-0.3
9	0.3	22	1.6	6.3	-0.5
10	0.4	21	1.4	5.9	-0.2

Case	X_1	X_2	X_3	X_4	X_5
1	1	0	0	0	0
2	0	1	0	0	0
3	0	0	0	0	1
4	0	0	0	1	0
5	0	0	1	0	0
6	0	0	0	0	0
7	0	0	0	0	0
8	0	0	0	0	0
9	0	0	0	0	0
10	0	0	0	0	0

6

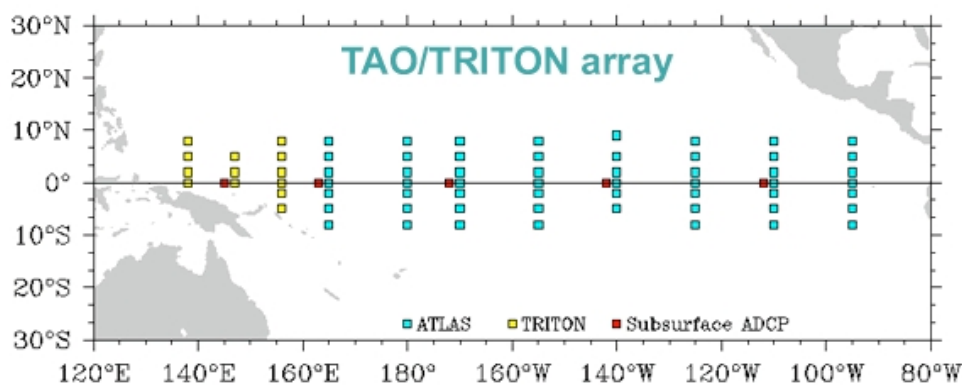
EXAMPLE

Tropical Atmosphere-Ocean Array

Number of cases: 736

Number of variables: 8

Sea Surface Temp, Air Temp,
Humidity, UWind, VWind + Year,
Lat Long



7

OVERVIEW

1993 El Nino

1997 Normal

Variable	Number of missing values	
	1993	1997
sea surface temp	3	0
air temp	4	77
humidity	93	0
uwind	0	0
vwind	0	0

R package: norm

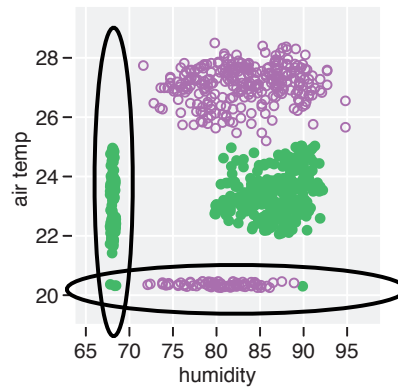
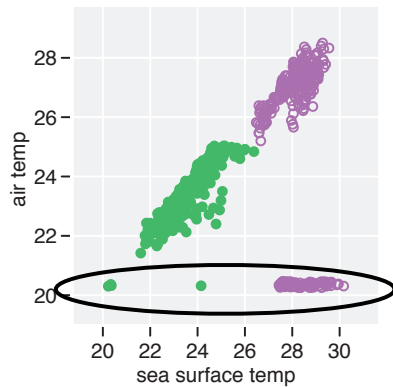
No. of missings on a case	1993		1997	
	No. of cases	%	No. of cases	%
3	2	0.5	0	0
2	2	0.5	0	0
1	90	24.5	77	20.9
0	274	74.5	291	79.1

8

USING THE MARGINS

1993 El Nino

1997 Normal

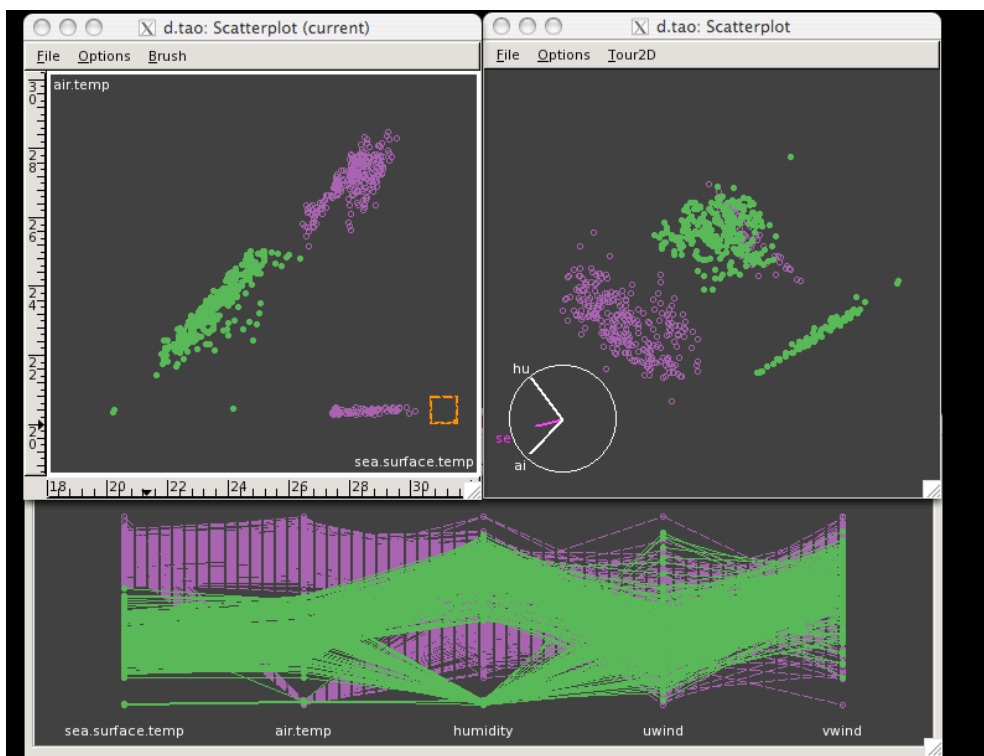


Association between temperatures. Years separated. More missings on _____ than _____.

Missings on _____ only occur in 1997.

9

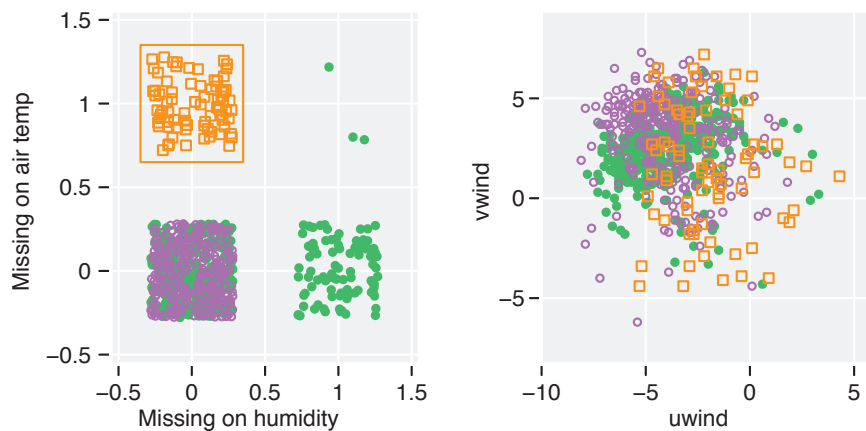
LIMITATION



Missings look like _____ in high-d plots, and in parallel coordinates they look like _____ at the very bottom.

10

TRACKING MISSING USING THE SHADOW MATRIX



Missings on air temp have _____
values on uwind than non-
missings.

11

MISSING STRUCTURE

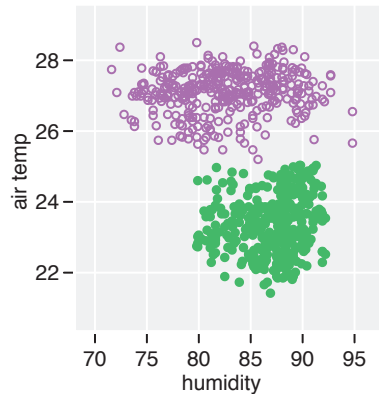
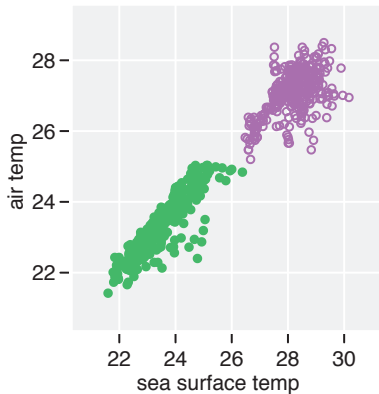
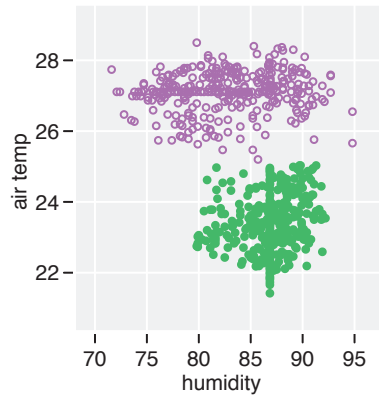
Missing values are _____!

Imputation will need to use
_____ of missing and not
missing.

12

IMPUTING MISSINGS

_____ for each year

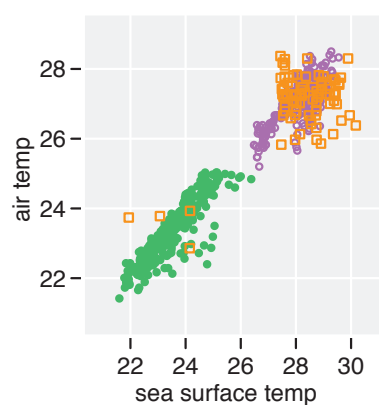
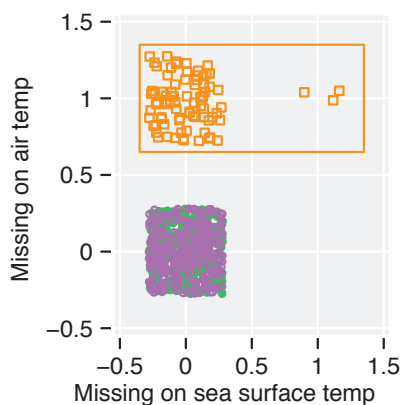


_____ from each year

What do you notice?

13

USING THE SHADOW MATRIX

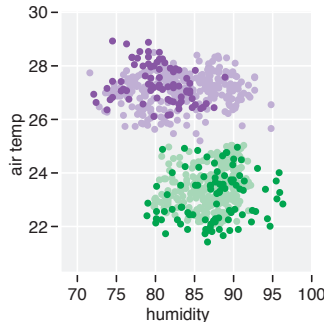
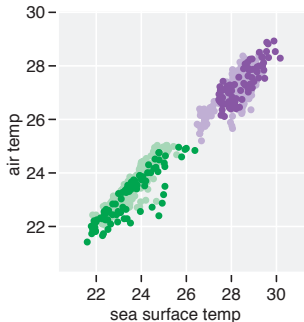


Imputed values which disappeared can be revealed by _____ on the shadow matrix.

14

MULTIPLE IMPUTATION

1



Missing values are imputed by simulating from a _____, having mean vector and variance-covariance matrix equal to the sample quantities. Sampling _____ times allows for estimating statistics for the missing values.

2

3, 4, 5, ...

15

SUMMARY

- _____ missings: by variable, by case
- _____ plots of missings, in the margins
- _____ summary statistics using as much data as possible.
- Determine _____ of missings: MAR, MCAR, MNAR
- Decide on a good way to _____ missings, as simple as possible with out affecting results.

16

This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 United States License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.