

# Data visualization

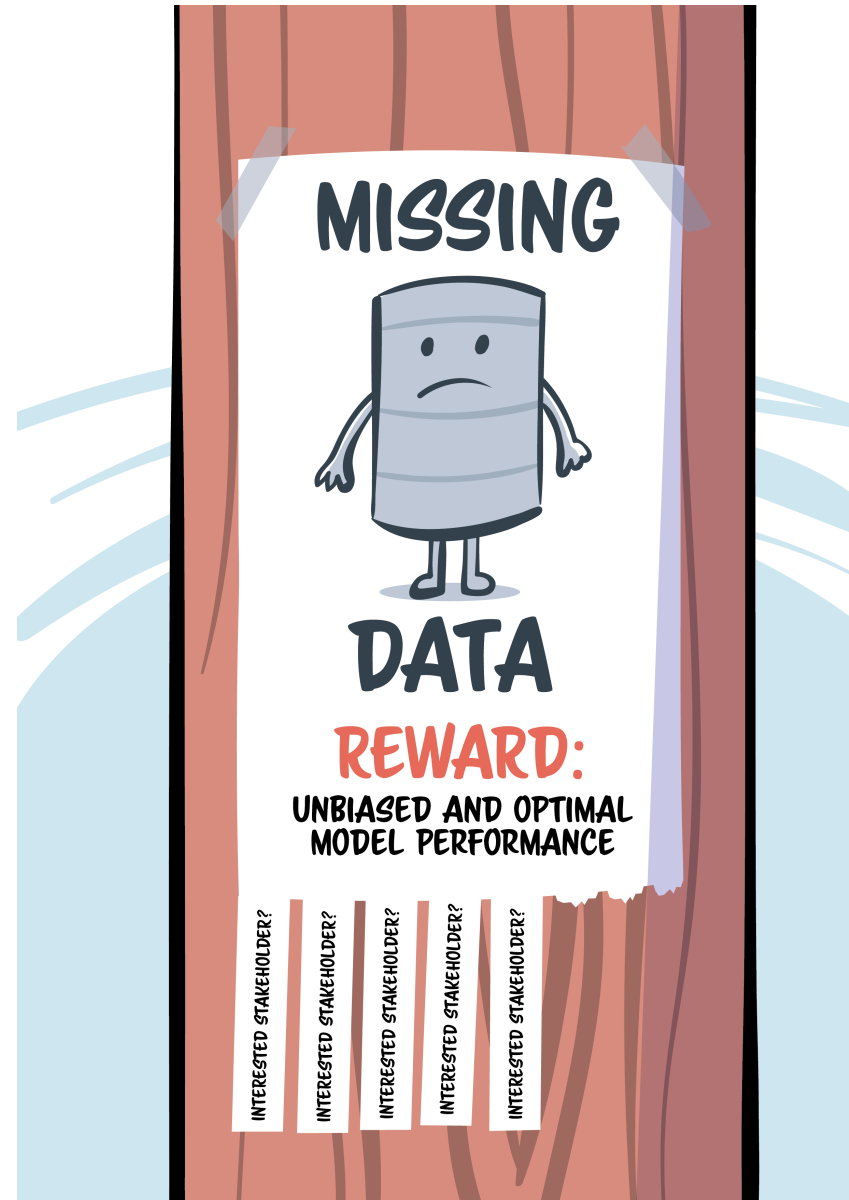
for incomplete datasets in R

Hanne Oberman

PhD candidate at Utrecht University

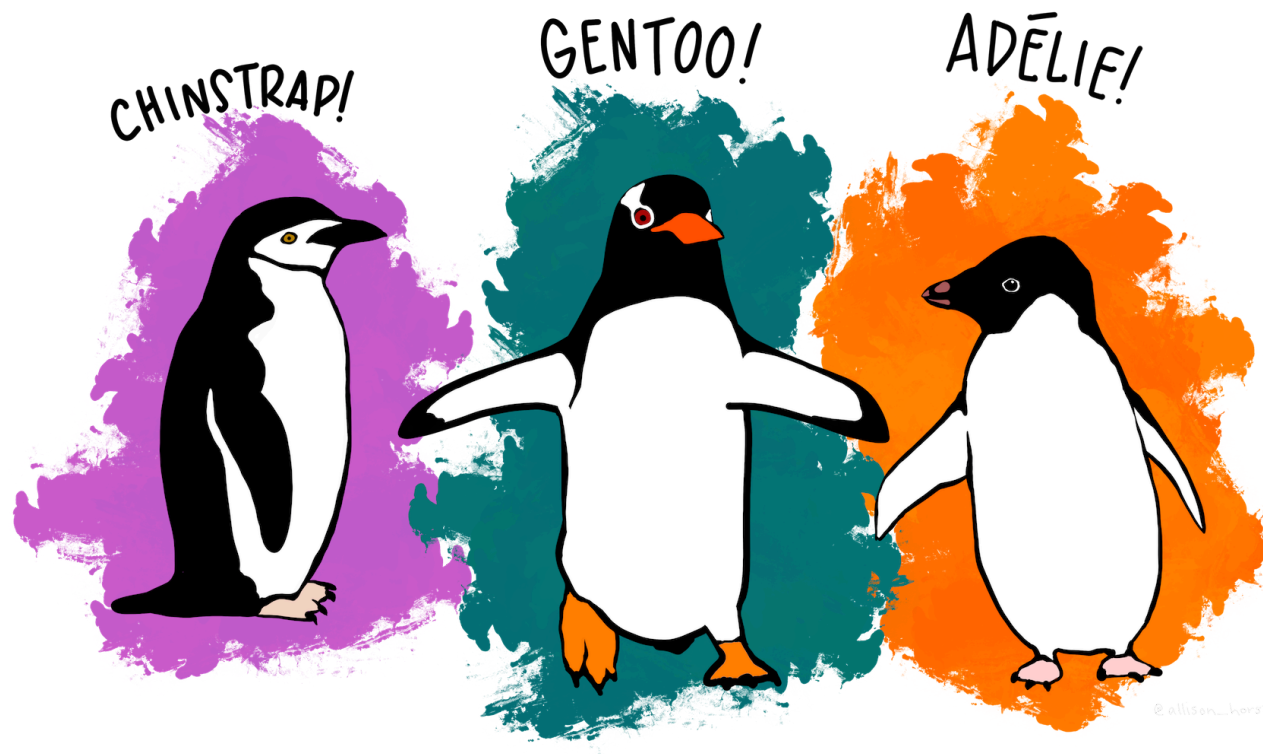


# Missingness



# Case study

```
1 set.seed(123)
2 library(palmerpenguins)
3 library(mice)
4 library(ggmice)
5 library(ggplot2)
```



# Incomplete data

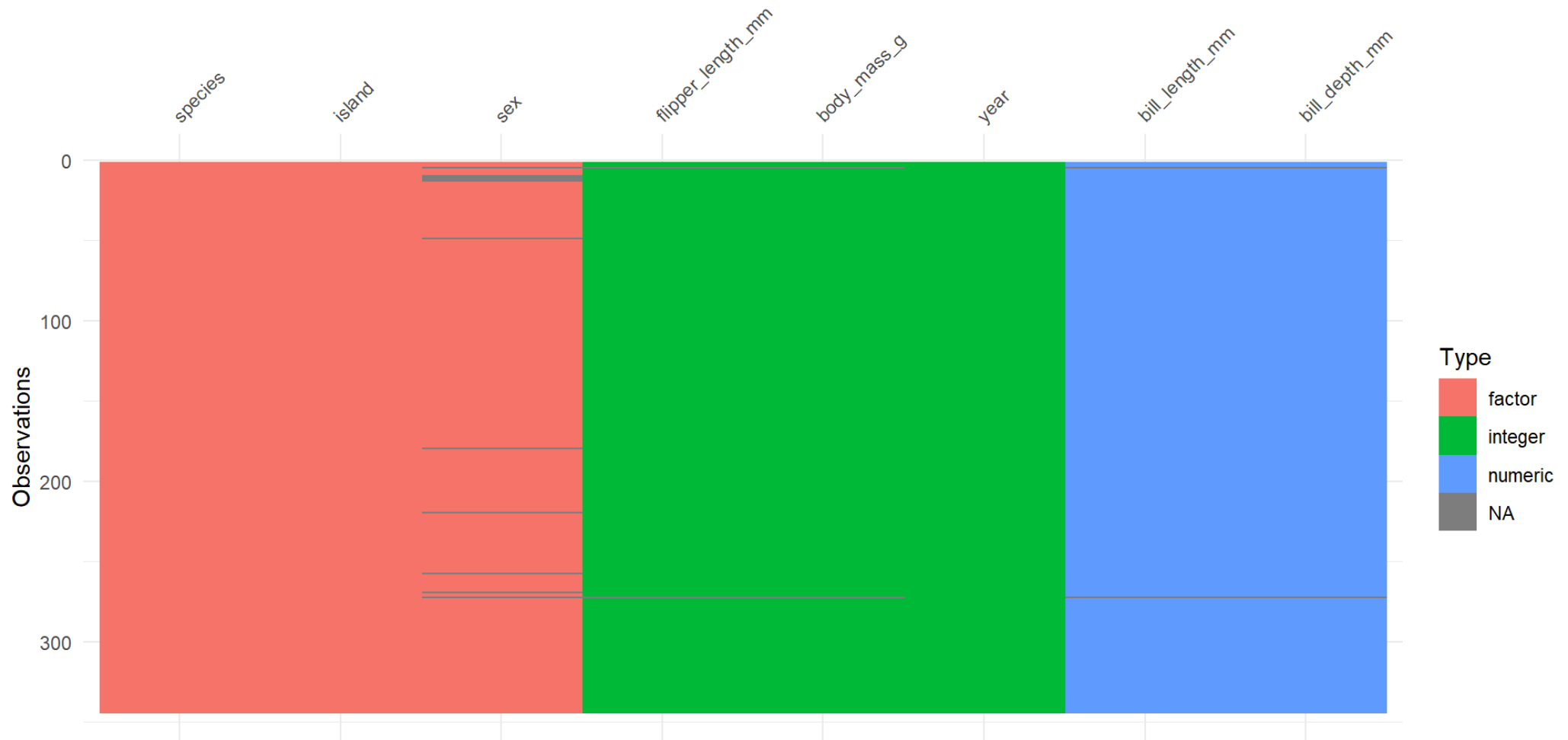
```
1 str(penguins)
```

```
tibble [344 × 8] (S3: tbl_df/tbl/data.frame)
 $ species      : Factor w/ 3 levels "Adelie","Chinstrap",...: 1 1 1 1 1 1
1 1 1 1 ...
 $ island       : Factor w/ 3 levels "Biscoe","Dream",...: 3 3 3 3 3 3 3 3
3 3 ...
 $ bill_length_mm : num [1:344] 39.1 39.5 40.3 NA 36.7 39.3 38.9 39.2 34.1
42 ...
 $ bill_depth_mm  : num [1:344] 18.7 17.4 18 NA 19.3 20.6 17.8 19.6 18.1
20.2 ...
 $ flipper_length_mm: int [1:344] 181 186 195 NA 193 190 181 195 193 190 ...
 $ body_mass_g    : int [1:344] 3750 3800 3250 NA 3450 3650 3625 4675 3475
4250 ...
 $ sex           : Factor w/ 2 levels "female","male": 2 1 1 NA 1 2 1 2 NA
NA ...
 $ year          : int [1:344] 2007 2007 2007 2007 2007 2007 2007 2007 2007
2007
```



# Incomplete data

```
1 visdat::vis_dat(penguins)
```



# Response indicator

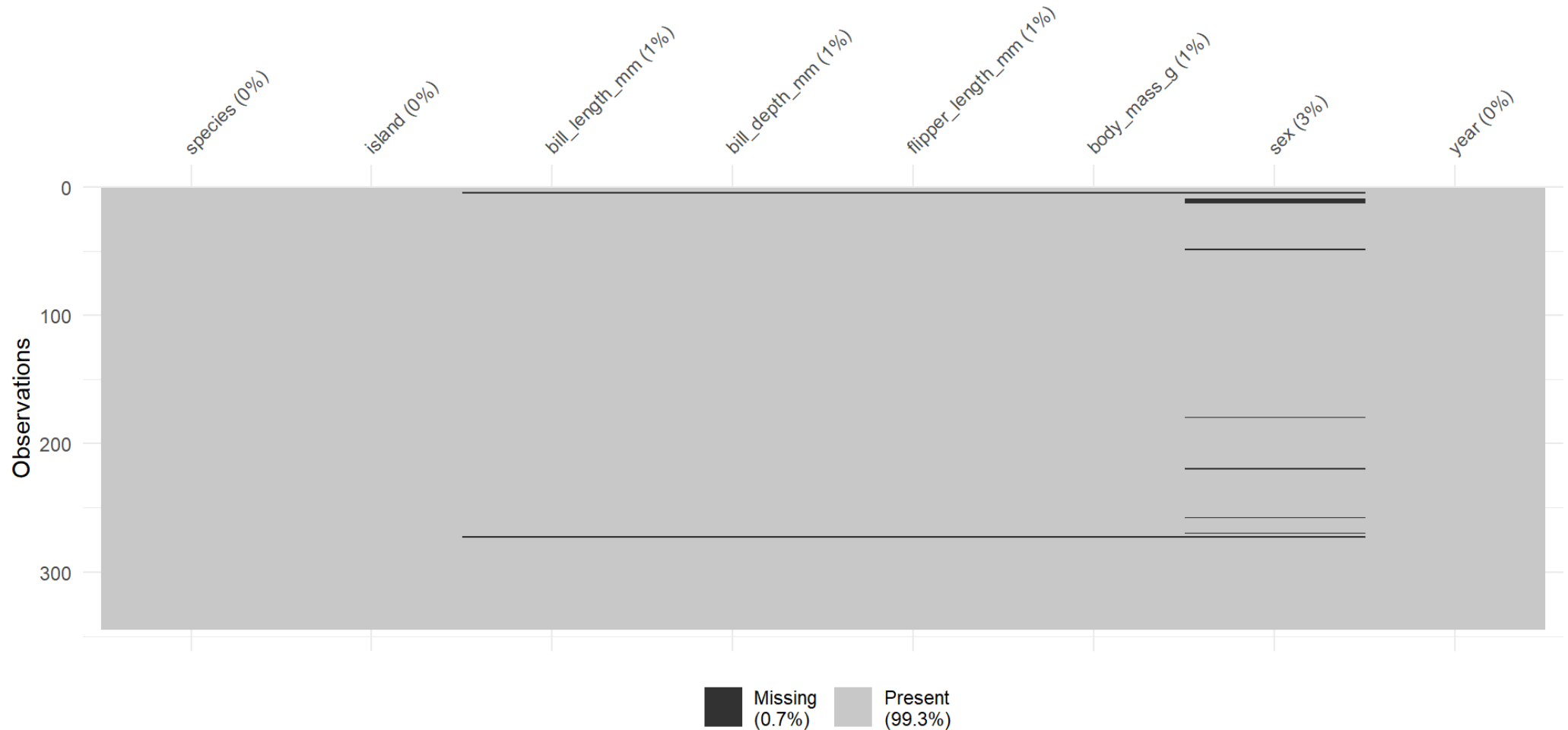
```
1 is.na(penguins)
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm
[1,]	FALSE	FALSE	FALSE	FALSE	FALSE
[2,]	FALSE	FALSE	FALSE	FALSE	FALSE
[3,]	FALSE	FALSE	FALSE	FALSE	FALSE
[4,]	FALSE	FALSE	TRUE	TRUE	TRUE
[5,]	FALSE	FALSE	FALSE	FALSE	FALSE
[6,]	FALSE	FALSE	FALSE	FALSE	FALSE
[7,]	FALSE	FALSE	FALSE	FALSE	FALSE
[8,]	FALSE	FALSE	FALSE	FALSE	FALSE
[9,]	FALSE	FALSE	FALSE	FALSE	FALSE
[10,]	FALSE	FALSE	FALSE	FALSE	FALSE
[11,]	FALSE	FALSE	FALSE	FALSE	FALSE
[12,]	FALSE	FALSE	FALSE	FALSE	FALSE
[13,]	FALSE	FALSE	FALSE	FALSE	FALSE
[14,]	FALSE	FALSE	FALSE	FALSE	FALSE
[15,]	FALSE	FALSE	FALSE	FALSE	FALSE



# Response indicator

```
1 naniar::vis_miss(penguins)
```



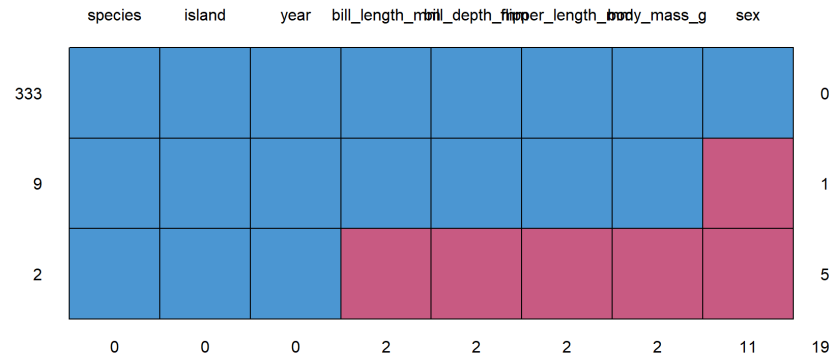
# Missing data pattern

```
1 md.pattern(penguins)
```

	species	island	year	bill_length_mm	bill_depth_mm	flipper_length_mm
333	1	1	1	1	1	1
9	1	1	1	1	1	1
2	1	1	1	0	0	0
	0	0	0	2	2	2

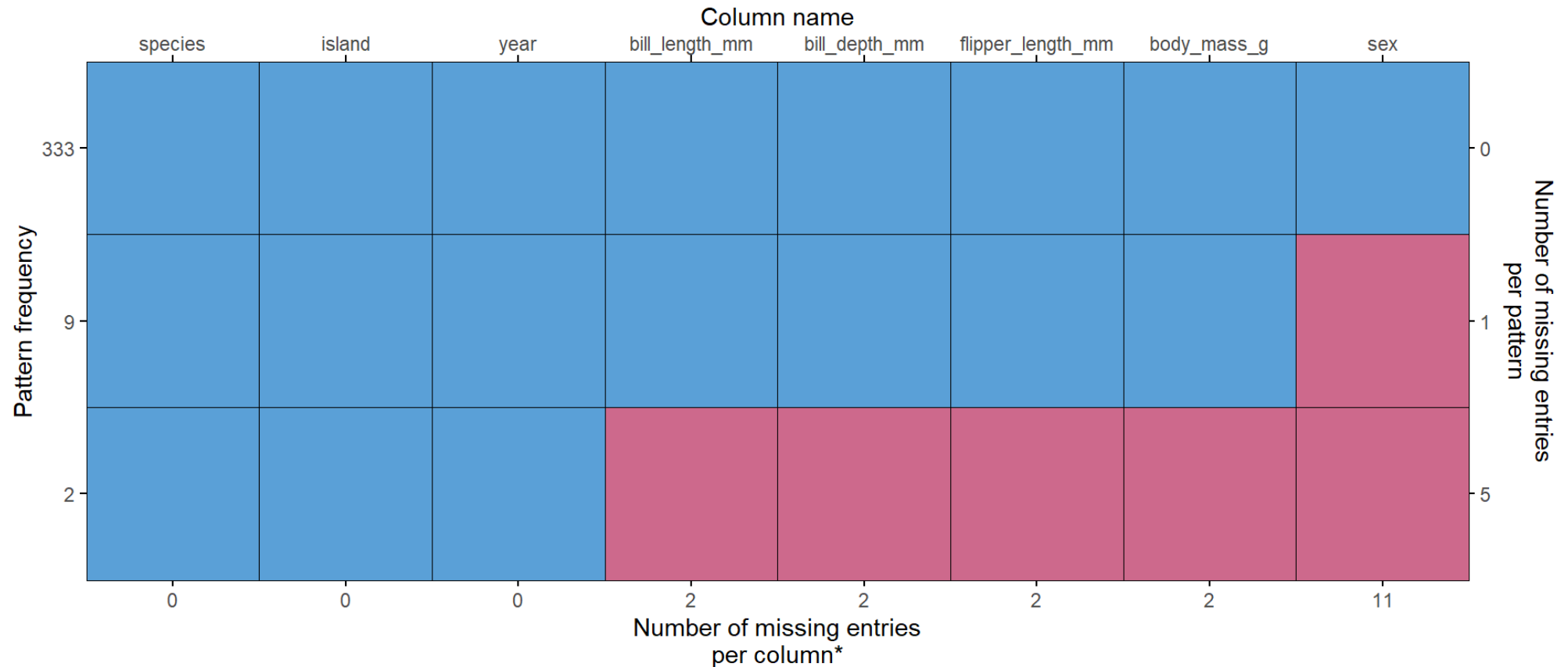
	body_mass_g	sex
333	1	0
9	1	1
2	0	5
	2	11 19





# Missing data pattern

```
1 plot_pattern(penguins)
```

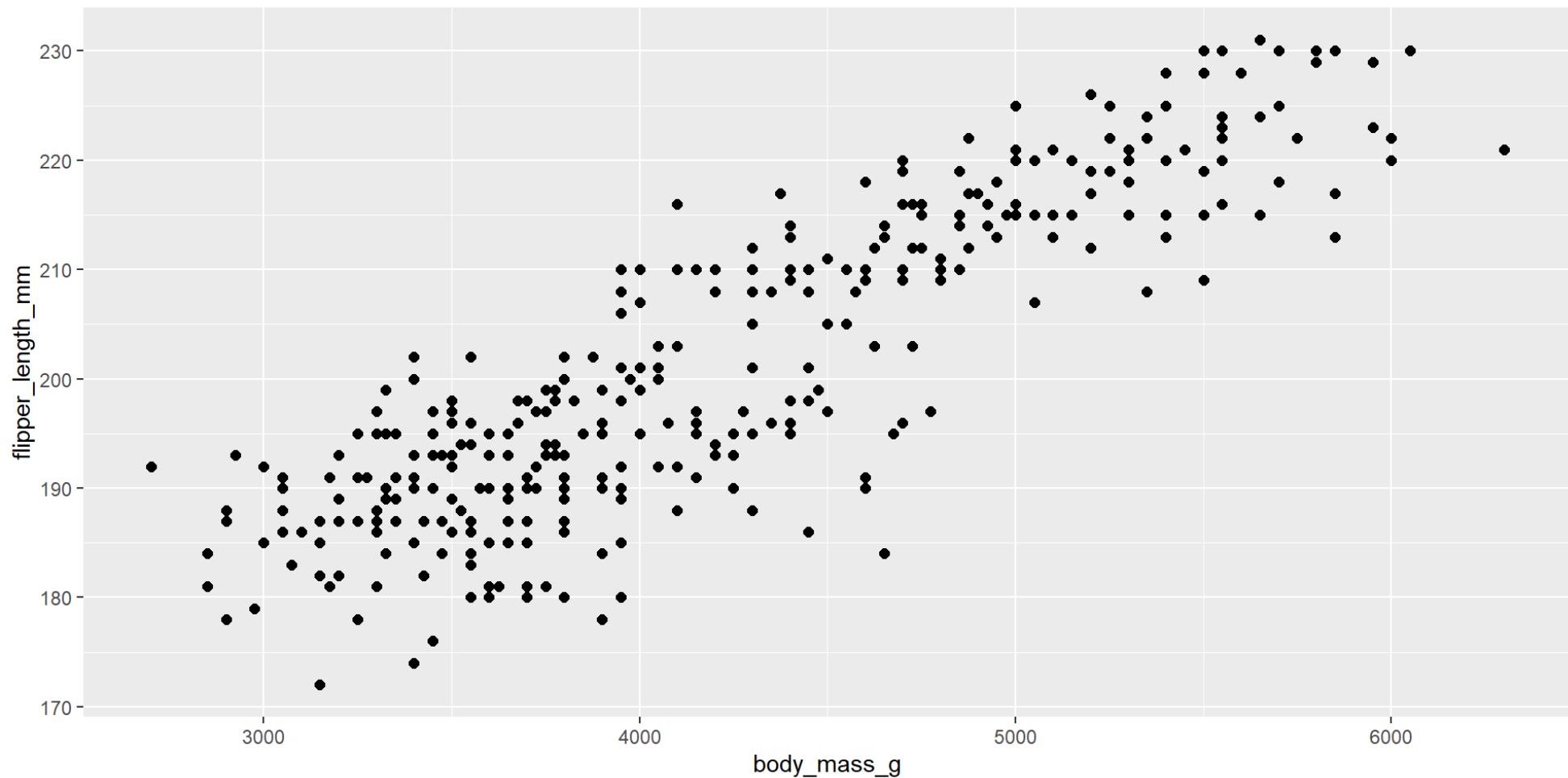


\*total number of missing entries: 19



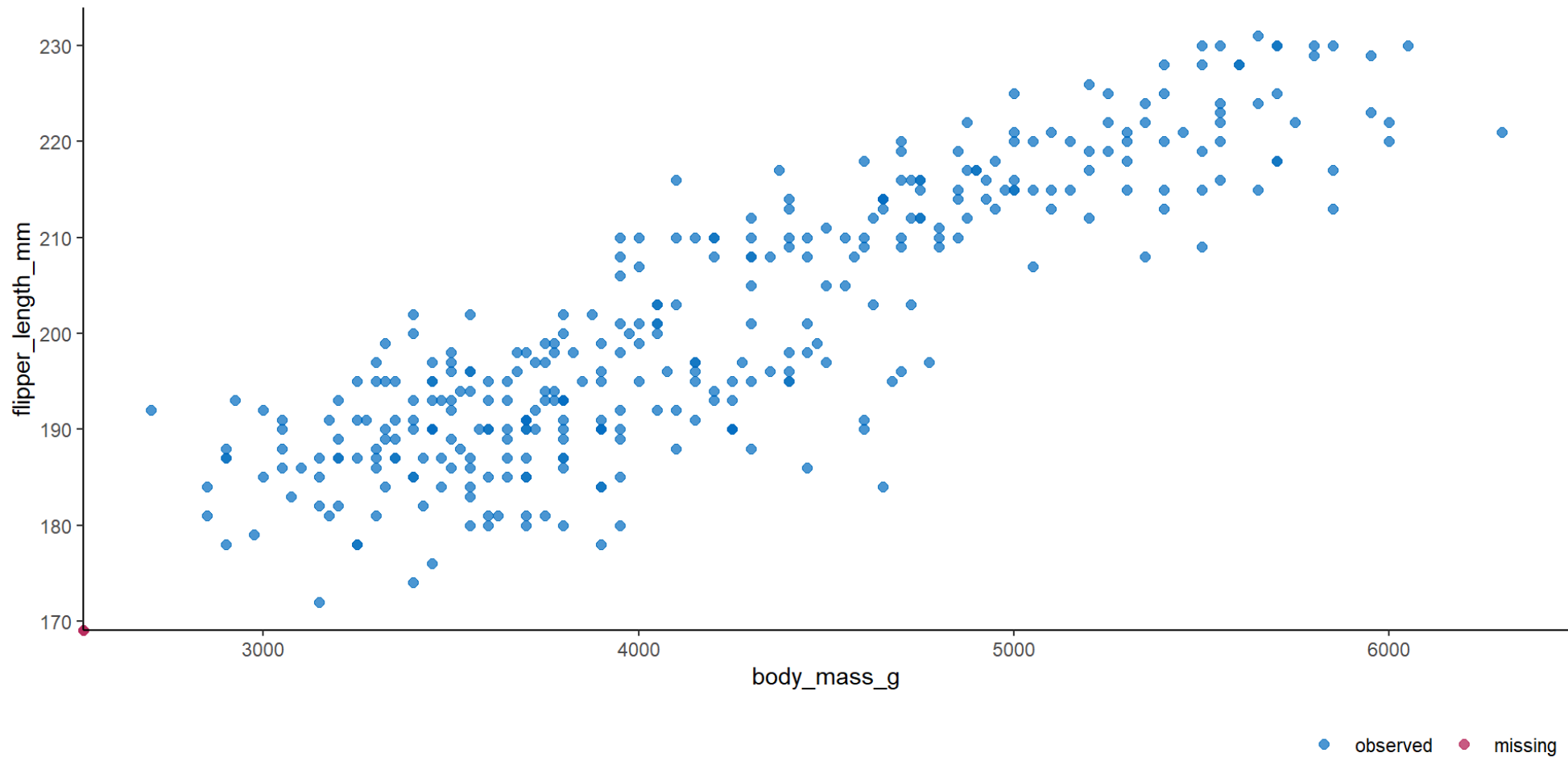
# Scatter plot

```
1 ggplot(penguins, aes(body_mass_g, flipper_length_mm)) +  
2   geom_point(size = 2)
```



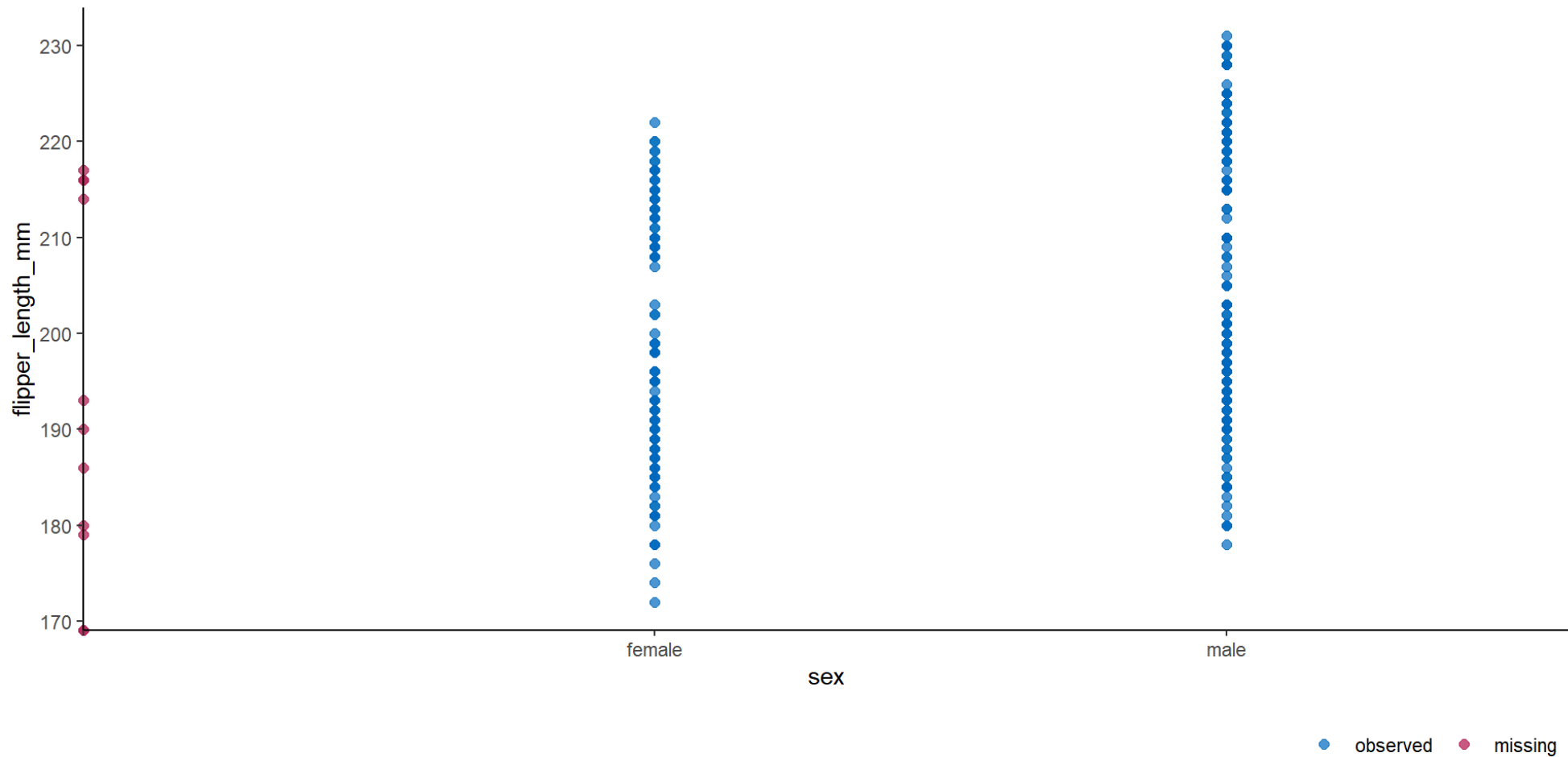
# Scatter plot

```
1 ggmlce(penguins, aes(body_mass_g, flipper_length_mm)) +  
2   geom_point(size = 2)
```



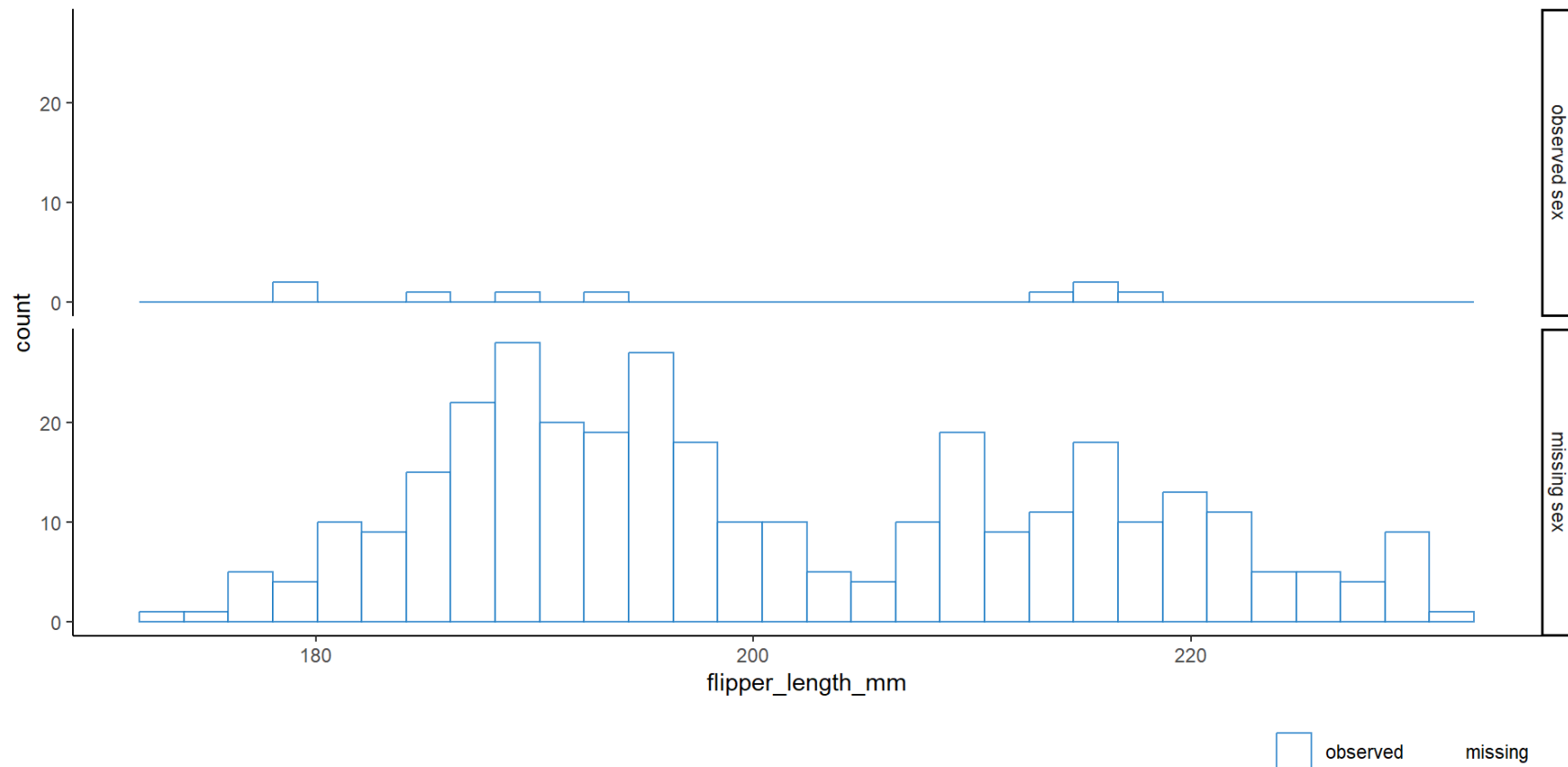
# Scatter plot

```
1 ggmlce(penguins, aes(sex, flipper_length_mm)) +  
2   geom_point(size = 2)
```



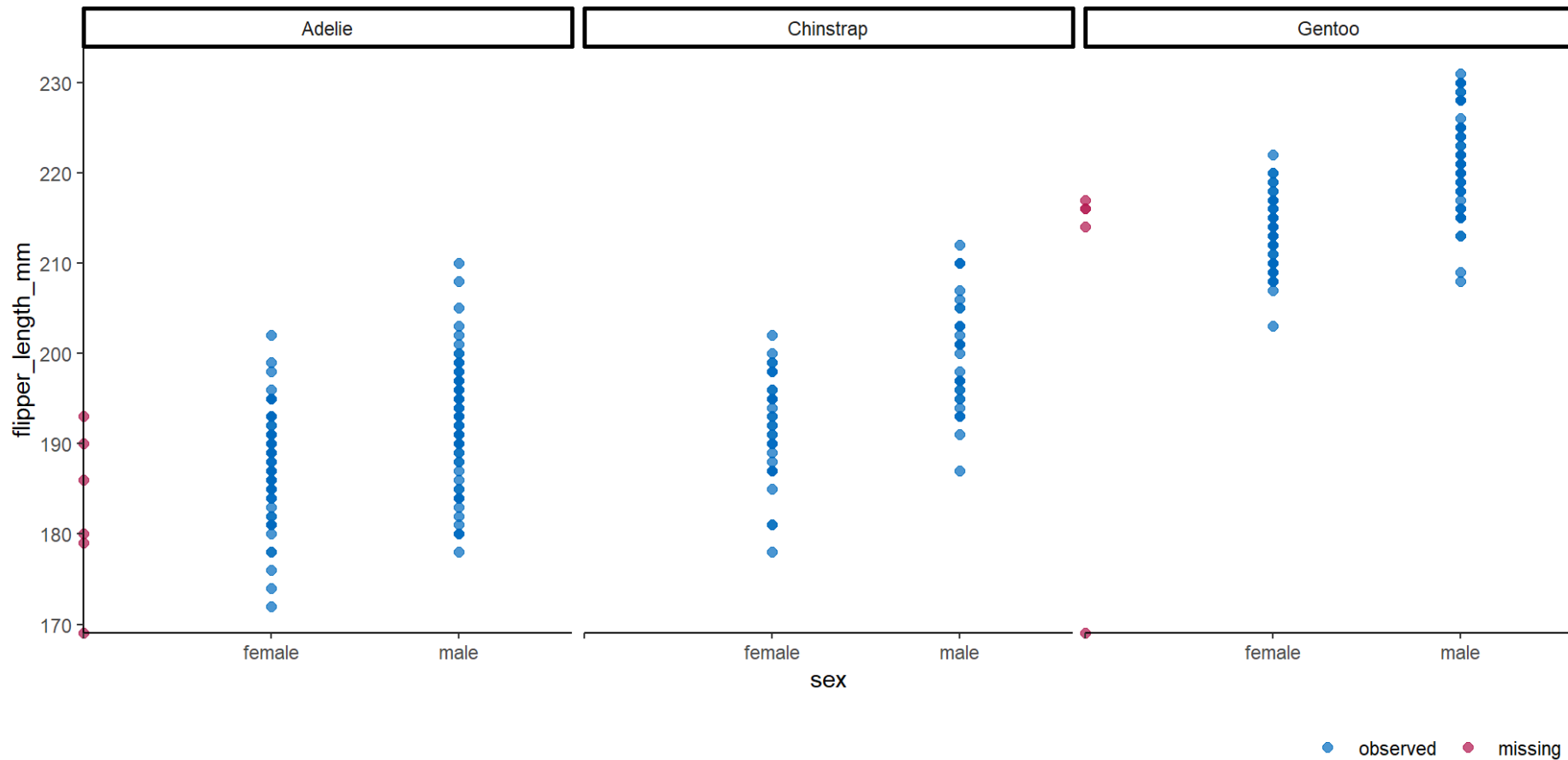
# Faceted distribution

```
1 ggmgice(penguins, aes(flipper_length_mm)) +  
2   geom_histogram(fill = "white") +  
3   facet_grid(  
4     factor(is.na(sex) == 0, labels = c("observed sex", "missing sex")) ~ .)
```



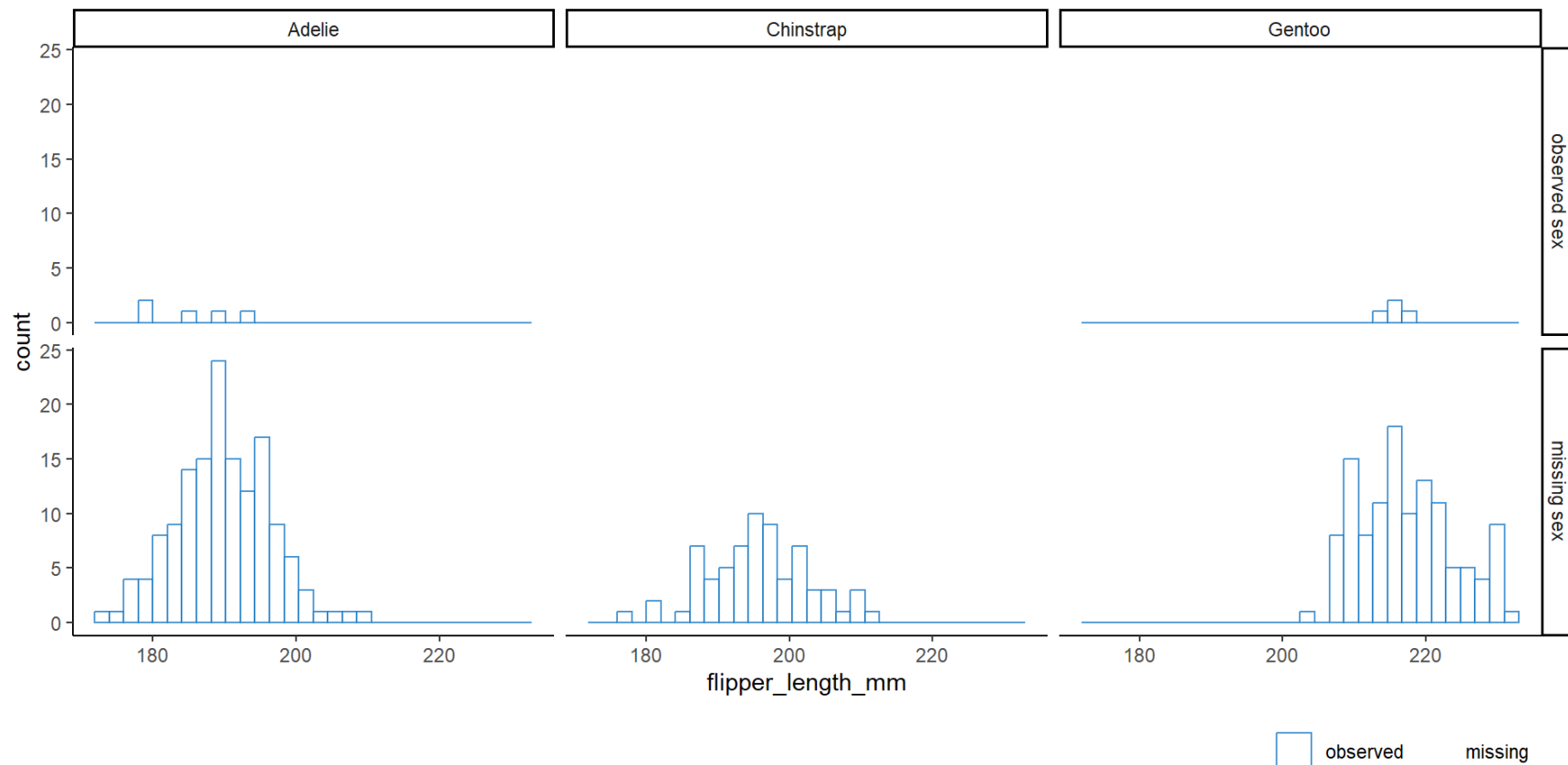
# Faceted scatter plot

```
1 ggmlce(penguins, aes(sex, flipper_length_mm)) +  
2   geom_point(size = 2) +  
3   facet_wrap(~species)
```



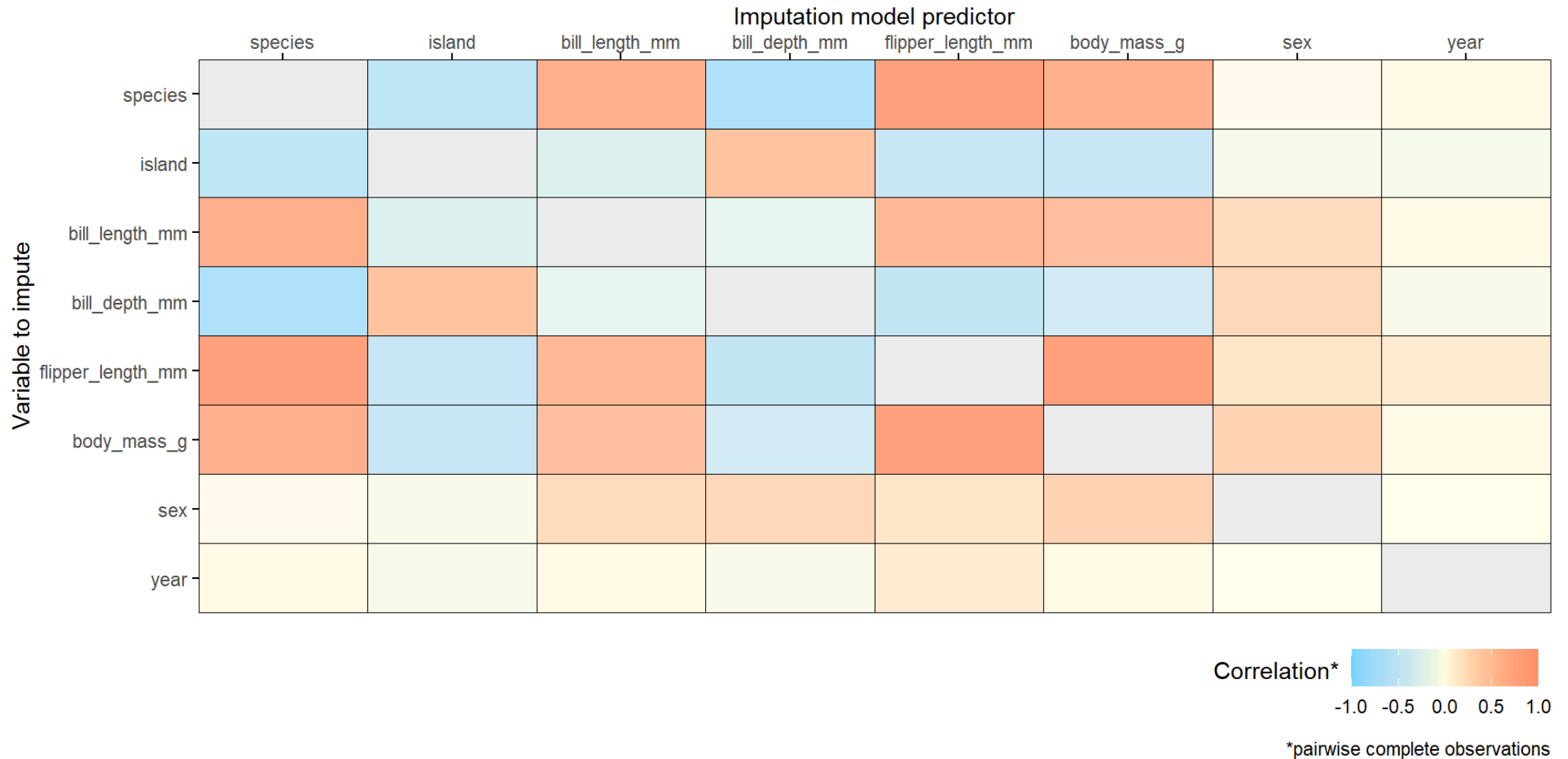
# Faceted distribution

```
1 ggmgice(penguins, aes(flipper_length_mm)) +  
2   geom_histogram(fill = "white") +  
3   facet_grid(  
4     factor(is.na(sex) == 0, labels = c("observed sex", "missing sex")) ~ sp
```



# Correlation

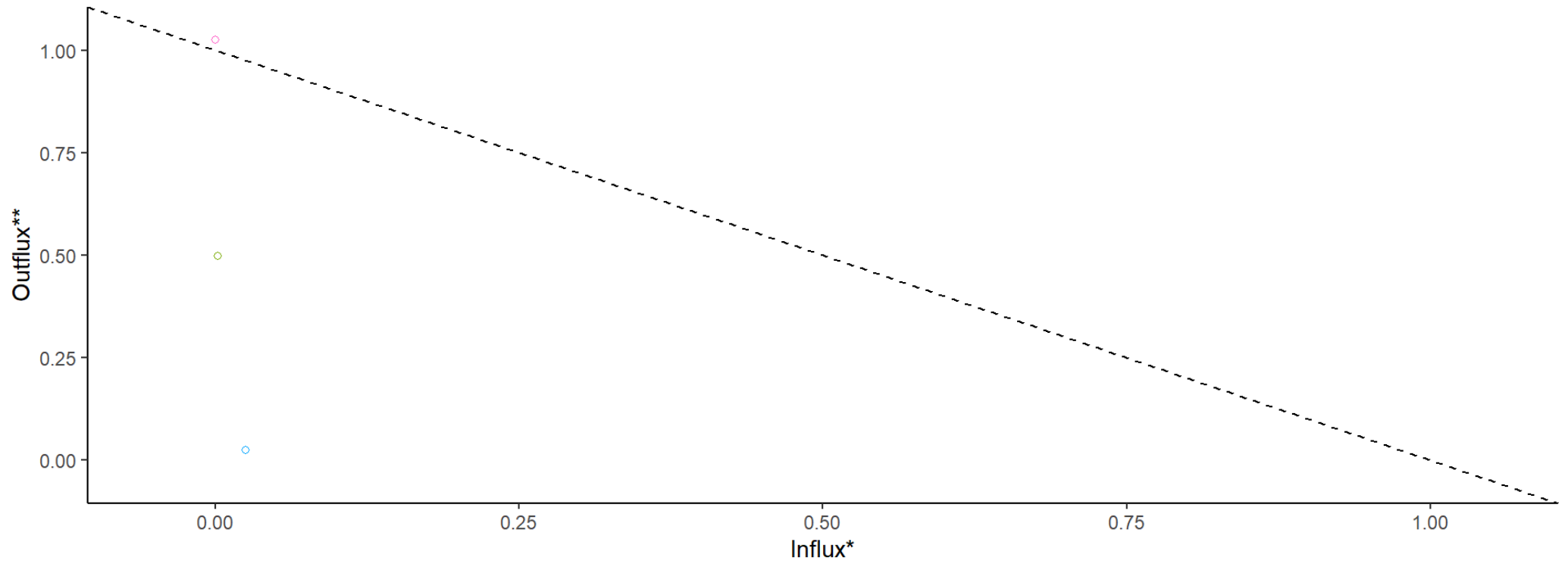
```
1 plot_corr(penguins, square = FALSE)
```





# Flux plot

```
1 plot_flux(penguins, label = FALSE)
```



○ bill\_depth\_mm    ○ body\_mass\_g    ○ island    ○ species  
○ bill\_length\_mm    ○ flipper\_length\_mm    ○ sex    ○ year

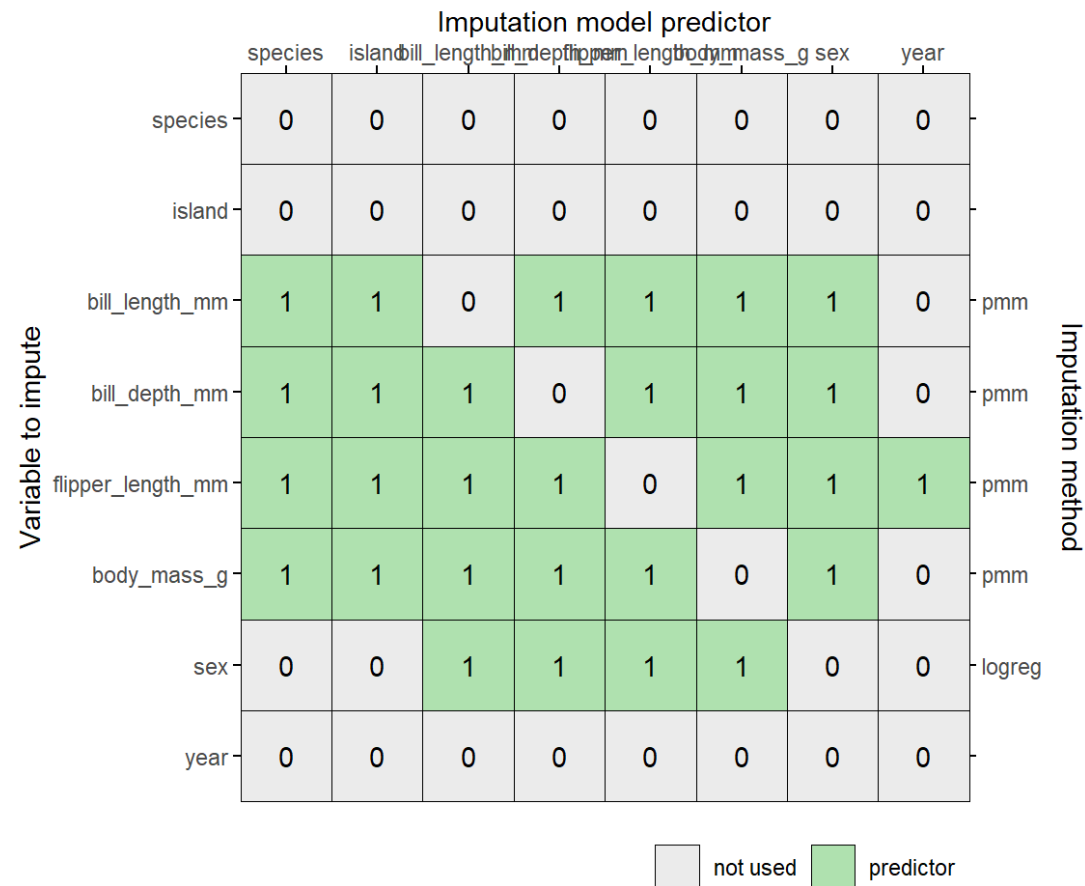
\*connection of a variable's missingness indicator with observed data on other variables

\*\*connection of a variable's observed data with missing data on other variables



# Imputation models

```
1 pred <- quickpred(penguins)
2 meth <- make.method(penguins)
3 plot_pred(pred, method = meth)
```



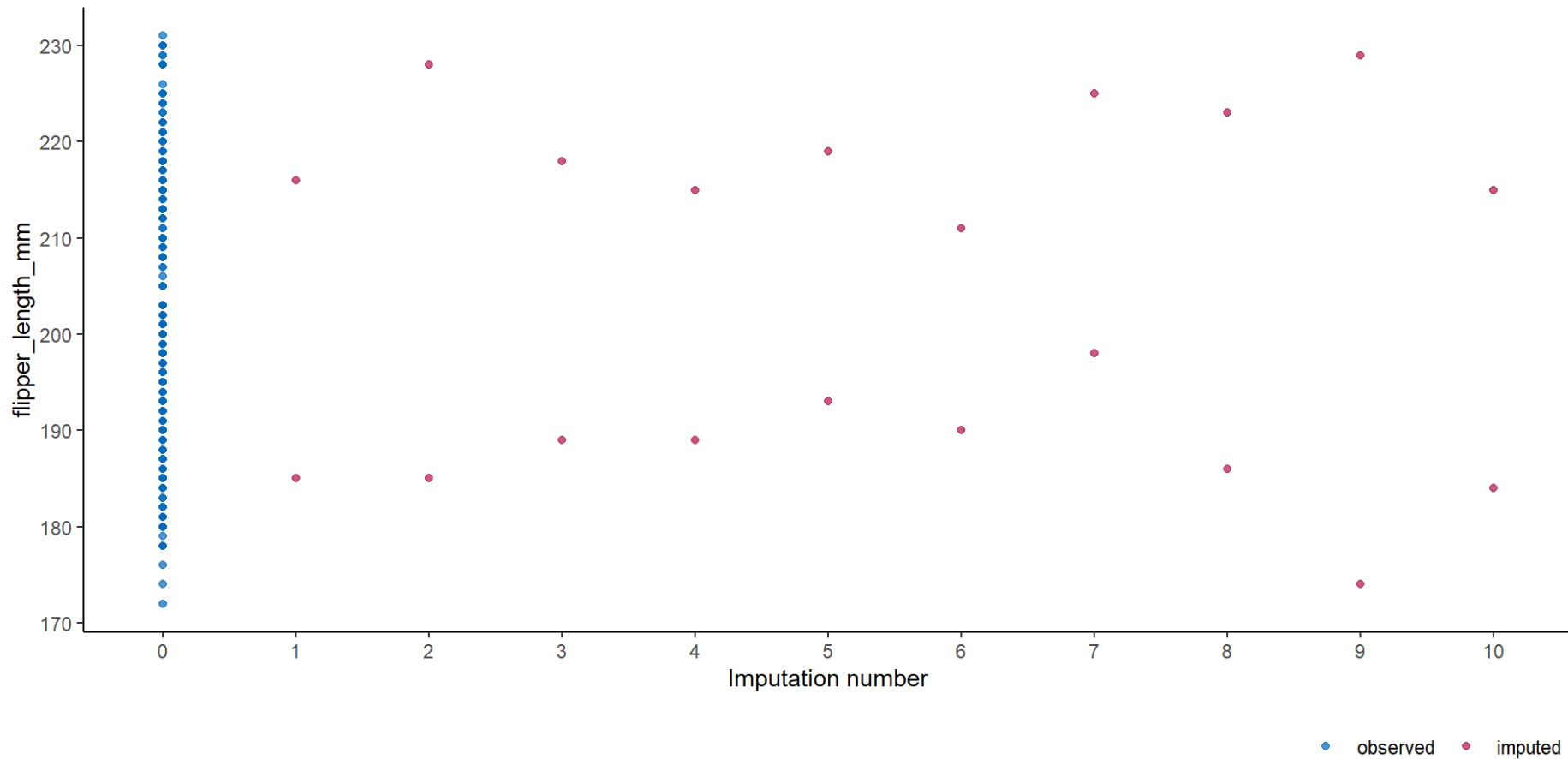
# Impute

```
1 imp <- mice(penguins, pred = pred, method = meth, m = 10, print = FALSE)
2 plot_trace(imp)
```



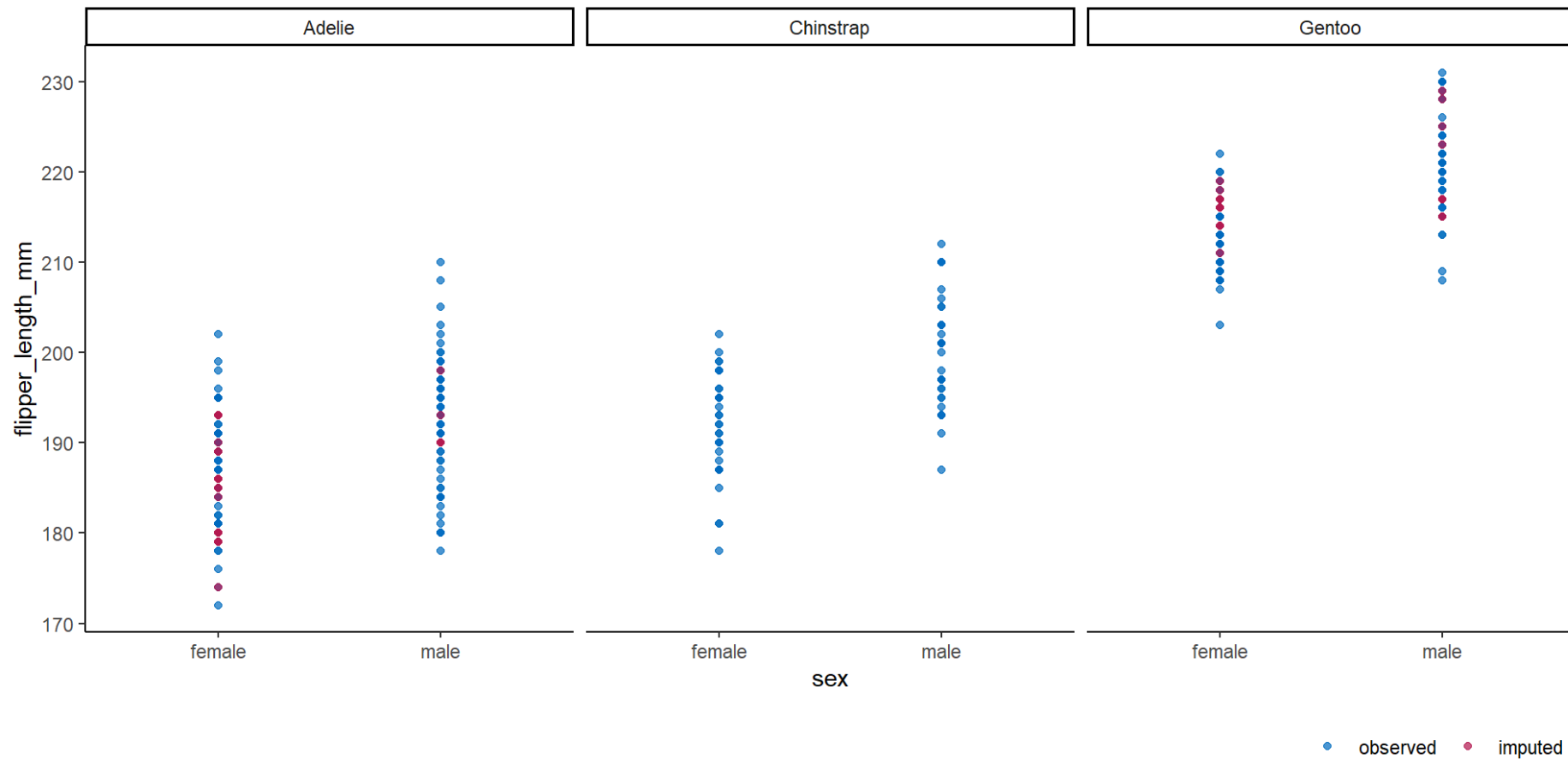
# Boxplot

```
1 ggmlce(imp, aes(x = .imp, y = flipper_length_mm)) +  
2   geom_point() +  
3   labs(x = "Imputation number")
```



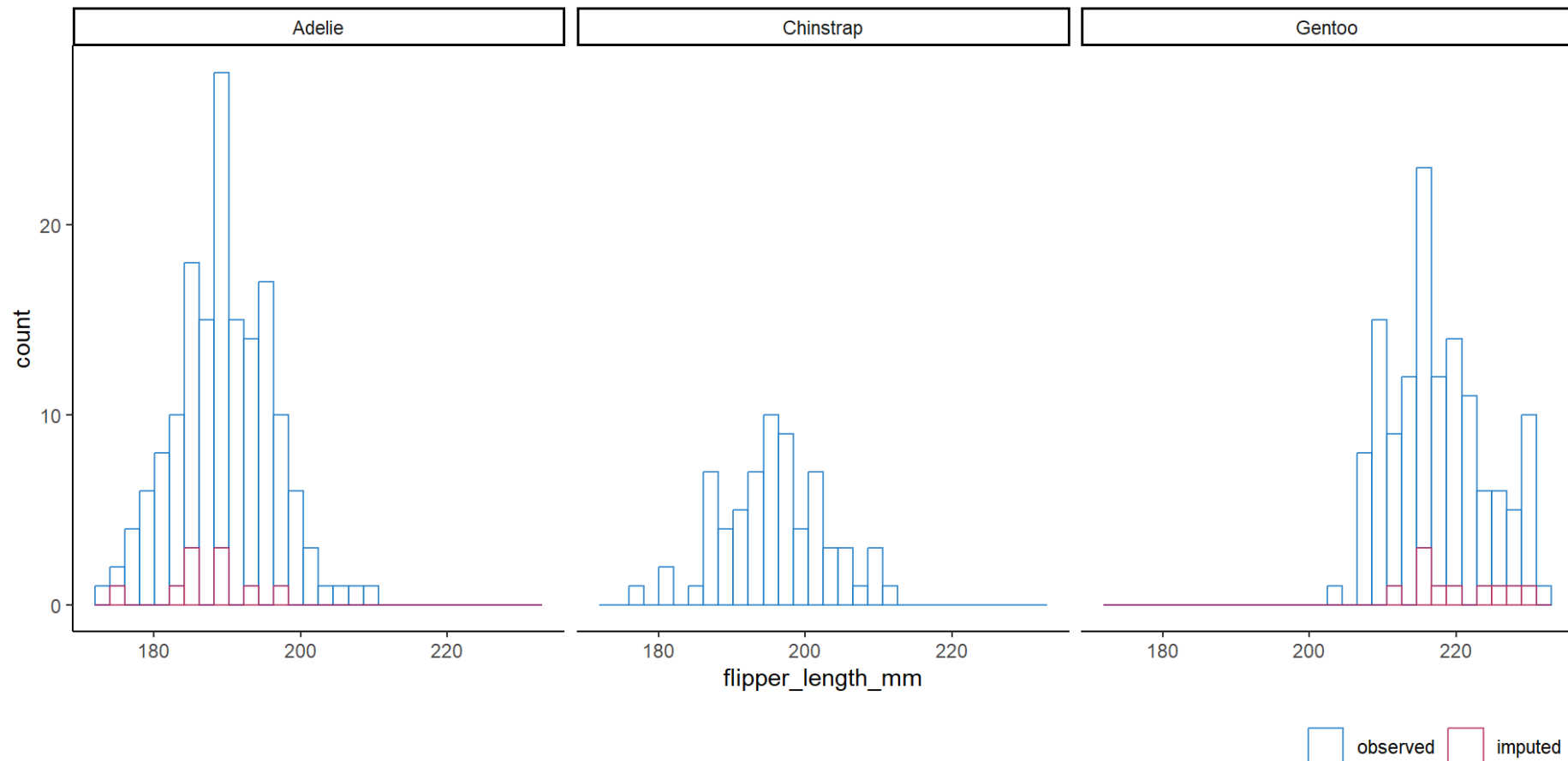
# Scatter plot

```
1 ggmlce(imp, aes(sex, flipper_length_mm)) +  
2   geom_point() +  
3   facet_grid(~species)
```



# Faceted distribution

```
1 ggmlce(imp, aes(flipper_length_mm)) +  
2   geom_histogram(fill = "white") +  
3   facet_grid(~ species)
```



# Thank you!



