


Data Science for Actuaries (ACT6100)

Arthur Charpentier

Non-Supervisé # 4 (k plus proches voisins & imputation)

automne 2Q20

 <https://github.com/freakonometrics/ACT6100/>

Missing Values & k -NN

There are 3 major types of missingness to be concerned about:

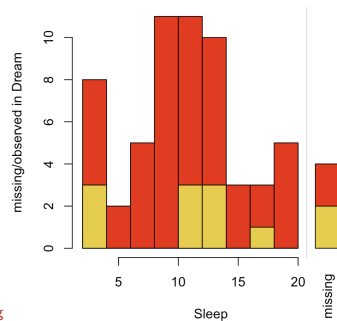
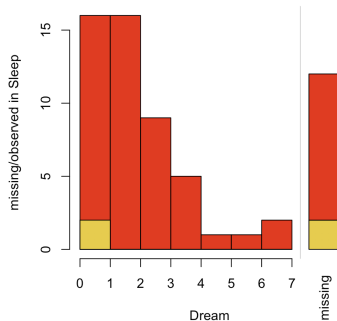
- ▶ **Missing Completely at Random** (MCAR) - the probability of missingness in a variable is the same for all units. Like randomly poking holes in a data set.
- ▶ **Missing at Random** (MAR) - the probability of missingness in a variable depends only on available information (in other predictors).
- ▶ **Missing Not at Random** (MNAR) - the probability of missingness depends on information that has not been recorded and this information also predicts the missing values.

Missing Values & k-NN

via Allison & Chichetti (1976)

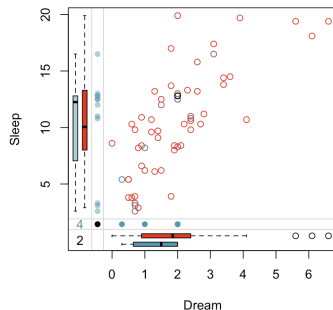
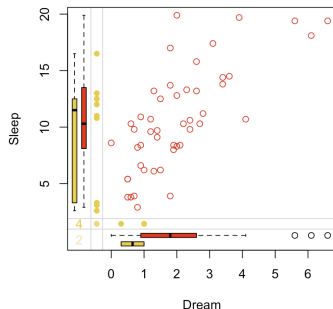
```
1 > library(VIM)
2 > x = sleep[, c("Dream", "
    Sleep")]
3 > summary(x)
4      Dream      Sleep
5  Min.    :0.000   Min.    : 2.60
6  1st Qu.:0.900   1st Qu.: 8.05
7  Median :1.800   Median :10.45
8  Mean   :1.972   Mean    :10.53
9  3rd Qu.:2.550   3rd Qu.:13.20
10 Max.    :6.600   Max.    :19.90
11 NA's    :12      NA's    :4
```

```
1 > histMiss(x)
2 > histMiss(x[,2:1])
```



Missing Values & k -NN

```
1 > marginplot(x)
2 > x_imputed = kNN(x)
3 > marginplot(x_imputed,
4   delimiter = "_imp")
5 > i=apply(x,1,function(x) sum(
6   is.na(x)))
7 > cor(x[i==0,])
8
9   Dream      Sleep
10  Dream 1.000000 0.727087
11  Sleep 0.727087 1.000000
12
13 > cor(x_imputed[,1:2])
14
15   Dream      Sleep
16  Dream 1.0000000 0.7396222
17  Sleep 0.7396222 1.0000000
```



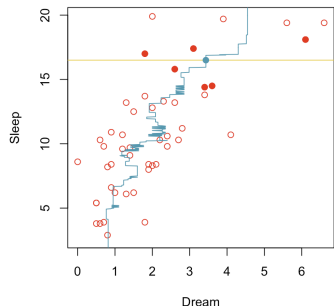
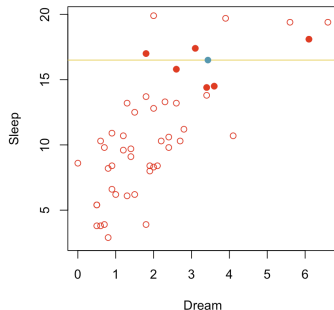
Missing Values & k -NN

If $x_{i,2}$ is missing

$$\bar{x}_{i,2} = \frac{1}{5} \sum_{j \in V} x_{j,2}$$

$$V = \{j : \|x_{j,1} - x_{i,1}\| \leq k\}$$

```
1 > i = which(is.na(x[,2]))[1]
2 > xc=x[!is.na(x[,2]),]
3 > R=rank(abs(xc[,1]-x[i,1]),
           ties.method = "random")
4 > ic = which(R<=5)
5 > mean(xc[ic,2])
```



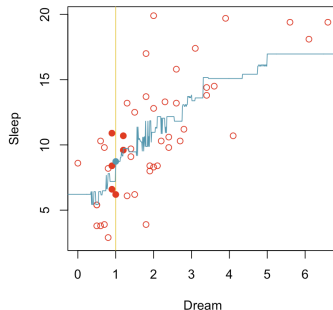
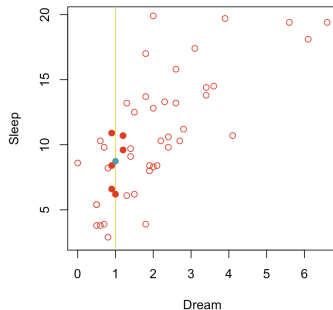
Missing Values & k -NN

If $x_{i,1}$ is missing

$$\bar{x}_{i,1} = \frac{1}{5} \sum_{j \in V} x_{j,1}$$

$$V = \{j : \|x_{j,2} - x_{i,2}\| \leq k\}$$

```
1 > i = which(is.na(x[,1]))[1]
2 > xc=x[!is.na(x[,1]),]
3 > R=rank(abs(xc[,2]-x[i,2]),
   ties.method = "random")
4 > ic = which(R<=5)
5 > mean(xc[ic,1])
```



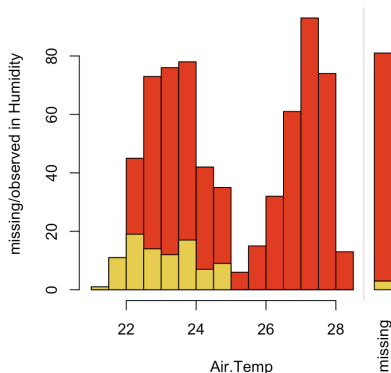
Missing Values & k -NN

```
1 > library(VIM)
2 > data(tao)
3 > y = tao[, c("Air.Temp", "Humidity")]
4 > summary(y)
5      Air.Temp      Humidity
6  Min.      :21.42   Min.      :71.60
7  1st Qu.:23.26   1st Qu.:81.30
8  Median :24.52   Median :85.20
9  Mean    :25.03   Mean     :84.43
10 3rd Qu.:27.08   3rd Qu.:88.10
11 Max.     :28.50   Max.      :94.80
12 NA's     :81     NA's      :93
```

Missing Values & k -NN

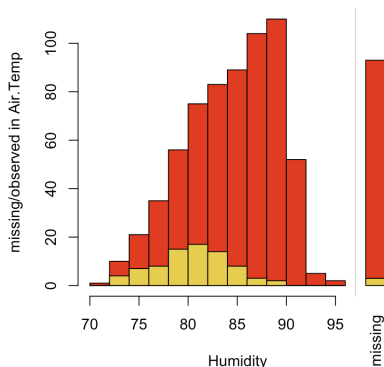
Missing humidity given
the temperature

```
1 > y = tao[,c("Air.Temp", "
    Humidity")]
2 > histMiss(y)
```



Missing temperature given
the humidity

```
1 > y = tao[,c("Humidity", "
    Air.Temp")]
2 > histMiss(y)
```



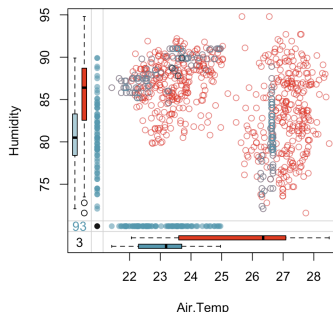
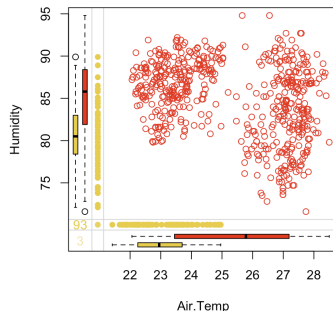
Missing Values & k -NN

This package contains a k -Nearest Neighbors algorithm for imputation

```
1 > tao_kNN = kNN(tao, k = 5)
```

Imputation can be visualized using

```
1 vars = c("Air.Temp", "Humidity",  
           ", "Air.Temp_imp", "  
           Humidity_imp")  
2 marginplot(tao_kNN[,vars],  
             delimiter="imp", alpha  
             =0.6)
```



Missing Values & k -NN

This package contains a k -Nearest Neighbors algorithm for imputation

```
1 > tao_kNN = kNN(tao, k = 5)
```

Imputation can be visualized using

```
1 vars = c("Air.Temp", "Humidity",  
           ", "Air.Temp_imp", "  
           Humidity_imp")  
2 marginplot(kNN(tao[vars[1:2]],  
              k = 5), delimiter="imp",  
              alpha=0.6)
```

