# Data Science for Actuaries (ACT6100)

Arthur Charpentier
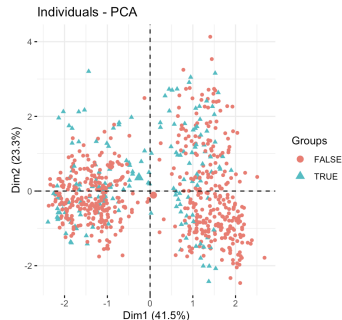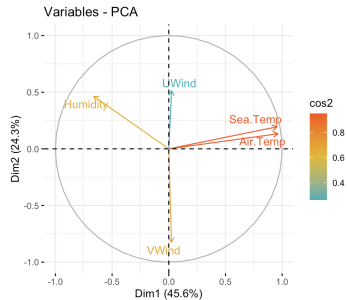
Non-Supervisé # 5 (ACP & imputation)

automne 2Q20

⬡ https://github.com/freakonometrics/ACT6100/

# Missing Values & PCA

```
1 > library(missMDA)
2 > library(VIM)
3 > names(tao)[4]="Sea.Temp"
4 > res.pca = PCA(tao[,4:8],graph=
      FALSE)
5 Warning message:
6 In PCA(tao[, 4:8], graph = FALSE)
7   Missing values are imputed by
      the mean of the variable: you
      should use the imputePCA
      function of the missMDA
      package
8 > fviz_pca_var(res.pca, col.var =
      "cos2")
9 > miss.ind = apply(tao[,4:8],1,
      function(x) sum(is.na(x))>0)
10 > fviz_pca_ind(res.pca, label="
      none", habillage=miss.ind)
```



Variables - PCA



Individuals - PCA

## Missing Values & PCA

The goal of PCA is maximize dispersion (inertia) of projections, or equivalently to minimize distance between observations, and their projections : we approximate our dataset $\boldsymbol{X}$ ($n \times p$) with some lower rank matrix,

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n \times k}} \left\{ \|\boldsymbol{X} - \boldsymbol{Z}\|^2 \right\}$$
subject to rank($\boldsymbol{Z}$) $\leq s$

for some $s < p$, where $\|\boldsymbol{M}\| = \text{trace}(\boldsymbol{M}\boldsymbol{M}^\top)$
From singular value decomposition,

$$\boldsymbol{Z}_s^\star = \boldsymbol{U}_{n \times s} \boldsymbol{\Delta}_{s \times s} \boldsymbol{V}_{p \times s}^\top = \underbrace{\boldsymbol{F}_{n \times s}}_{=\text{PC scores}} \underbrace{\boldsymbol{V}_{p \times s}^\top}_{\text{loadings}}$$

That was possible only with complete information (no missing data)

# Missing Values & PCA

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n \times k}} \left\{ \|\boldsymbol{X} - \boldsymbol{Z}\|^2 \right\}$$
subject to rank$(\boldsymbol{Z}) \leq s$

becomes, with missing data, $\boldsymbol{W} = (\boldsymbol{W}_{ij})$, $\boldsymbol{W}_{ij} = \mathbf{1}(\boldsymbol{X}_{ij}$ missing$)$,

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n \times k}} \left\{ \|\boldsymbol{W} \odot (\boldsymbol{X} - \boldsymbol{Z})\|^2 \right\}$$
subject to rank$(\boldsymbol{Z}) \leq s$

This can be solved by iterative PCA, see Kiers (1997)
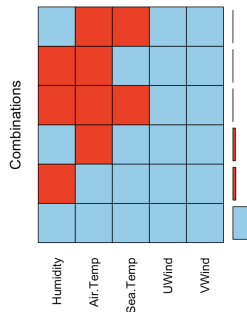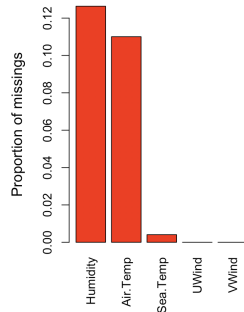
---

**Algorithm 1:** Iterative PCA

---

1 initialization : $\boldsymbol{X}^{(0)}$ completed by mean imputation, $s < p$;
2 **for** $t=1,2,...$ **do**
3    PCA on completed data, or SVD $\boldsymbol{U}_{n \times n}^{(t)}, \boldsymbol{\Delta}_{n \times p}^{(t)}, \boldsymbol{V}_{p \times p}^{(t)\top}$);
4    impute values with $Y^{(t)} = \boldsymbol{U}_{n \times s}, \boldsymbol{\Delta}_{s \times s}, \boldsymbol{V}_{p \times s}^{\top}$);
5    $\boldsymbol{X}^{(t)} \leftarrow \boldsymbol{W} \odot \boldsymbol{X} + (1 - \boldsymbol{W}) \odot Y^{(t)}$

---

# Missing Values & PCA



```
1  > res = summary(aggr(tao[,4:8],
       sortVar = TRUE),col=colrvim)
       $combinations
2
3  Variables sorted by number of
       missings:
4  Variable       Count
5  Humidity 0.126358696
6  Air.Temp 0.110054348
7  Sea.Temp 0.004076087
8     UWind 0.000000000
9     VWind 0.000000000
```
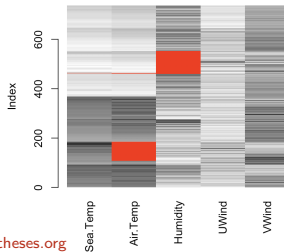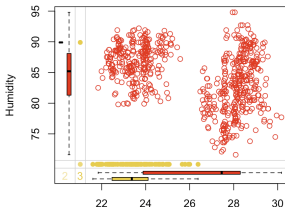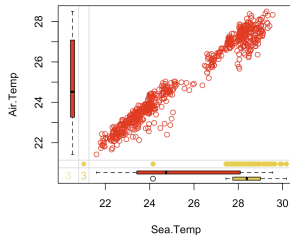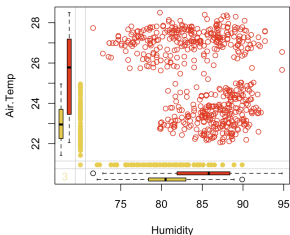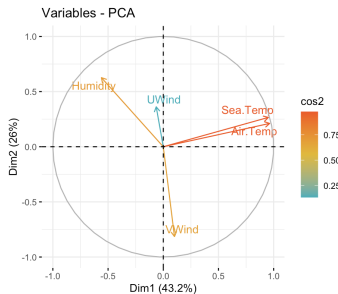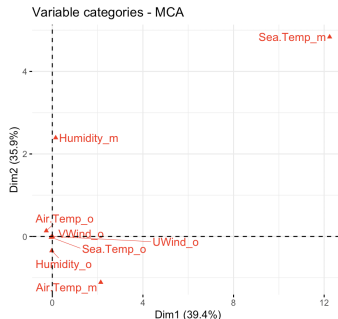
# Missing Values & PCA

```
1 > marginplot(tao[ ,c("Humidity", "Air.Temp")])
2 > marginplot(tao[ ,c("Sea.Temp", "Air.Temp")])
3 > marginplot(tao[ ,c("Sea.Temp", "Humidity")])
4 > matrixplot(tao[,4:8])
```
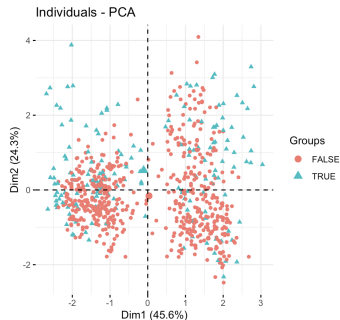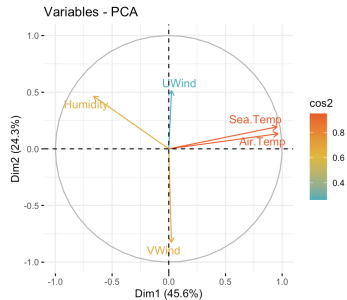
# Missing Values & PCA

```
1  > mis.ind = matrix("o", nrow =
      nrow(tao[,4:8]), ncol = ncol(
      tao[,4:8]))
2  > mis.ind[is.na(tao[,4:8])] = "m"
3  > dimnames(mis.ind) = dimnames(
      tao[,4:8])
4  > library(FactoMineR)
5  > res.mca = MCA(mis.ind)
6  > fviz_mca_var(res.mca, repel =
      TRUE)
7  > res.pca = PCA(tao[!miss.ind
      ,4:8],graph=FALSE)
8  > fviz_pca_var(res.pca, col.var =
      "cos2")
```



Variable categories - MCA



Variables - PCA

# Missing Values & PCA

```
1 > res.comp <- imputePCA(tao
     [,4:8], ncp = 2)
2 > res.comp$completeObs[1:3, ]
3      Sea.Temp Air.Temp Humidity
     UWind VWind
4 [1,]    27.59    27.15    79.6
     -6.4   5.4
5 [2,]    27.55    27.02    75.8
     -5.3   5.3
6 [3,]    27.57    27.00    76.5
     -5.1   4.5
7 >
8 > res.pca = PCA(res.
     comp$completeObs,graph=FALSE)
9 > fviz_pca_var(res.pca, col.var =
     "cos2")
10 > fviz_pca_ind(res.pca, label="
     none", habillage=miss.ind)
```



Variables - PCA



Individuals - PCA

# Missing Values & PCA

```
1 > tao_kNN = kNN(tao[,4:8], k = 5)
2 vars = c("Air.Temp","Humidity","
     Air.Temp_imp","Humidity_imp")
3 marginplot(tao_kNN[,vars],
     delimiter="_imp")
4 > X=data.frame(res.
     comp$completeObs)
5 names(X) = names(tao[,4:8])
6 W=matrix(FALSE, nrow = nrow(tao
     [,4:8]), ncol = ncol(tao
     [,4:8]))
7 W[is.na(tao[,4:8])] = TRUE
8 W=data.frame(W)
9 names(W) = paste(names(tao[,4:8])
     ,"_imp",sep="")
10 tao_ACP <- data.frame(X,W)
11 vars <- c("Air.Temp","Humidity","
     Air.Temp_imp","Humidity_imp")
12 marginplot(tao_ACP[,vars],
     delimiter="_imp"
```