

Fiche bilan SAE

Nom de la SAE	SAE Modèle linéaire		Semestre / Période	Semestre 4
Volume horaire consacré par l'étudiant	Avec enseignant	3h	En autonomie	3h
Coéquipiers :	Sami Said		Mohamed Belarbi	
	Franck Tankapanya			

Sujet spécifique	Expliquer ou prédire une variable quantitative à partir de plusieurs facteurs
Objectifs	<ul style="list-style-type: none"> • Importer et Comprendre les Données • Examiner les Relations entre Variables • Modéliser les Relations à travers la Régression Linéaire Multiple • Optimiser et Évaluer le Modèle

Fiche bilan SAE

Livrables

```
options(repos = c(CRAN = "https://cran.rstudio.com/"))
install.packages("readxl")

## Installation du package dans 'C:/Users/Feky/AppData/Local/R/win-library/4.3/'
## (car 'lib' n'est pas spécifique)

## Le package 'readxl' a été décompressé et les sommes MD5 ont été vérifiées avec succès
##
## Les packages binaires téléchargés sont dans
## C:/Users/Feky/AppData/Local/Temp/rtmp56Gb/downloaded_packages

library(readxl)

## Warning: le package 'readxl' a été compilé avec la version R 4.3.3

foot<-read_excel("C:/FBAHCK/taff/football_2009_2016.xlsx" )
data=subset(foot, select=c(Goals,Club))
```

Proposer un modele goals en fonction du club

```
mod_foot <- lm(Goals~Club ~ 1, data = data)
summary(mod_foot)

##
## Call:
## lm(formula = Goals ~ Club, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.893 -1.431 -0.204  0.420  4.587
##
## Coefficients:
##              (Intercept)              Club
##              1.2644              1.8916663 ***
## ---
## Signif. codes:  0 '0.001 ***' 0.01 '**' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.421 on 1999 degrees of freedom
## Adjusted R-squared:  0.0000000
## F-statistic: 0.0000000 on 1 DF, 0.0000000 probability <= 0.0000000
```

Estimation des espérances de Y

```
summary(mod_foot)$coefficients[,1]
```

##	ClubAC Ajaccio	ClubAC Milan
##	1.2643678	1.8916663
##	ClubAlaves	ClubAlmeria
##	1.6400000	1.4571629
##	ClubAngers	ClubArsenal
##	1.3620090	0.6562500
##	ClubAston Villa	ClubAthletic Bilbao
##	2.4977778	1.4285714
##	ClubAtletico Madrid	ClubBayer Leverkusen
##	1.3125000	2.0245898
##	ClubBirmingham	ClubBayern Munich
##	2.4607843	1.4060606
##	ClubBlackburn	ClubBolton Wanderers
##	1.6794872	4.0926829
##	ClubCardiff	ClubCharlton Athletic
##	1.1538462	0.0000000
##	ClubCeltic	ClubChelsea
##	2.4375000	3.3730570
##	ClubDerby County	ClubDoncaster Rovers
##	1.4509804	1.4204545

Fiche bilan SAE

Tester au risque 5

```
## c -> fit(mod_1, data = data)
anova1, mod_1
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	21551	218700.0	NA	NA	NA	NA
2	21389	244835.5	162	73824.52	39.86481	0

2 rows

Au moins une équipe a marqué plus que 0 but : dans notre cadre cette question n'est pas intéressante car on l'observe lors des estimations : question 2) b) On teste si les moyennes de chaque groupe est nulle **on a le même modèle mais sous contrainte**

```
mod_c2 <- lm(Goals~C(Club, base = 1), data)
summary(mod_c2)
```

```
##
##
## Call:
## lm(formula = Goals ~ C(Club, base = 1), data = data)
##
## Residuals:
##      [1] 12.00000      [2] 12.00000      [3] 12.00000
##      [4] 12.00000      [5] 12.00000      [6] 12.00000
##      [7] 12.00000      [8] 12.00000      [9] 12.00000
##      [10] 12.00000     [11] 12.00000     [12] 12.00000
##
##
```

On a mis alpha égale a 0, donc le Club de l'AC Milan est passé en Intercept. Ce qui a comparé les estimations des nombres de buts des différents clubs en la comparant a celle de l'AC Milan0. On observe donc que les clubs ayant une estimation positive ont marqué plus de buts en moyenne que L'AC Milan. Tandis que ceux qui ont une estimation négative ont marqué moins de buts en moyenne que l'AC Milan. Mais si l'on veut déduire que la différence de buts entre l'AC Milan et un des clubs est significative on doit avoir une p-valeur inférieure a 0,05. On voit donc que le FC Barcelone a marqué en moyenne 2,2 buts par match de plus que l'AC Milan et la p-valeur est de 5.41e-12, donc on rejette H0. On peut donc affirmer que la différence de buts marqués en moyenne entre le FC Barcelone et l'AC Milan est significative. Dans le cas inverse, si on s'intéresse au club de Cordoba on peut voir que l'estimation de buts moyens marqués par Cordoba est de 0.77 but de moins que l'AC Milan. Puis la p-valeur est de 0.03 donc on rejette H0, on peut donc affirmer que la différence de buts entre cordoba et l'AC Milan est significative.

la contrainte est alpha=0 Intercept nest plus la

Estimation des espérances de Y

summary(mod_c2)\$coefficients[,1]

```
##
##      (Intercept)
##      1.891562655
##      C(Club, base = 2)FC_Cordoba
##      -0.427198489
##      C(Club, base = 2)Atletico
##      0.731562655
##      C(Club, base = 2)Almeria
##      -0.434226488
##      C(Club, base = 2)Angers
##      -0.324975000
##      C(Club, base = 2)Archie-Angon
##      -1.291362655
##      C(Club, base = 2)Arenas
##      0.682131513
##      C(Club, base = 2)Atoton Villa
##      -0.402962655
##      C(Club, base = 2)Alcala
##      -0.579802655
##      C(Club, base = 2)Atletico Bilbao
##      0.132849599
##      C(Club, base = 2)Atletico Madrid
##      0.569218840
##      C(Club, base = 2)Augsburg
##      -0.405090559
##      C(Club, base = 2)Augsburg
##      -0.212879886
##      C(Club, base = 2)Augsburg
##      -0.212879886
##
```

Tester au risque 5

m1 <- lm(Goals~1, data)
anova(m1, mod_c2)

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	21550	251544.4	NA	NA	NA	NA
2	21389	244835.5	161	6708.914	3.640344	9.897197e-49

2 rows

Au moins un club a mis plus ou moins de buts par rapport au Milan AC avec une p valeur de 2.2e-16 on rejette H0. La pvalue extrêmement faible de 2.2e-16 obtenue lors de la comparaison entre le modèle avec seulement l'Intercept et le modèle incluant la variable 'Club' en tant que facteur suggère que l'ajout du facteur 'Club' apporte une contribution significative à l'explication de la variance dans les objectifs marqués. Ainsi, nous rejetons l'hypothèse nulle selon laquelle tous les clubs marquent plus ou moins le même nombre de buts par rapport au Milan AC, avec un risque d'erreur de 5%.

Bilan de la SAE

(reproduire le tableau autant de fois que de compétences mobilisées dans la SAÉ)

Compétence	Modéliser les données dans un cadre statistique
Apprentissages critiques sollicités	Comprendre l'impact du type de données sur le choix de la modélisation à mettre en œuvre
	Apprécier les limites de validité et les conditions d'application d'un modèle
	Réaliser l'importance de la mise en œuvre d'une procédure de test statistique pour valider ou non une hypothèse
Composantes essentielles à respecter	En utilisant le modèle de données adapté aux besoins

Fiche bilan SAE

	En maîtrisant la qualité du modèle
	En choisissant le modèle adapté à la situation

Ma démarche

Savoirs / connaissances	Savoir-faire	Savoir-être
Statistiques descriptives, corrélation, régression linéaire, sélection de variables.	<ul style="list-style-type: none">• Manipuler• Programmer• Interpréter	<ul style="list-style-type: none">• Travail d'équipe,• Persévérance dans L'analyse de données• Communication claire des résultats.

Evaluation du résultat

- Ce que je trouve bien réalisé, pourquoi ?

Fiche bilan SAE

La méthodologie complète de modélisation, depuis l'analyse initiale jusqu'à la sélection de variables, permet une compréhension approfondie des dynamiques salariales. La collaboration en équipe a également enrichi l'approche analytique.

- Ce que je n'ai pas bien compris ; ce qui serait à améliorer pour une prochaine fois : pourquoi ? comment ?

L'un des aspects qui pourrait ne pas être immédiatement clair est la manière dont les hypothèses sous-jacentes à l'ANOVA ont été vérifiées. L'ANOVA repose sur certaines hypothèses importantes comme l'homogénéité des variances et la normalité des distributions des groupes. Pour une prochaine fois, il serait bénéfique d'inclure des tests spécifiques (par exemple, le test de Levene pour l'homogénéité des variances) ou des visualisations (QQ-plots pour la normalité) pour valider ces hypothèses. Ceci est crucial pour s'assurer que les conclusions tirées du modèle sont valides.

Éléments de preuve, ce que je peux montrer (*Choisir des éléments précis à mettre annexe*)

- 1) Extrait de code

Fiche bilan SAE

```
install.packages("readxl")
library(readxl)
foot<-read_excel("D:/BUT/2eme annee/modele lineaire/football_2009_2016.xlsx" )
data=subset(foot, select=c(Goals,Club))

# proposer un modele goals en fonction du club les moyennes par groupe de nos variables ne sont pas nul
mod_foot <- lm(Goals~Club -1, data = data)
summary(mod_foot)

#1 la contrainte est b=0 intercept nest plus la

#estimation des espérance de Y
summary(mod_foot)$coefficients[,1]

# Tester au risque 5

m0 <- lm(Goals~0, data = data)
anova(m0, mod_foot)

# au moins une equipe a marque plus que 0 but ; dans notre cadre cette question était pas intéressante car
#b=0 on teste si les moyenne de chaque groupe est nulle

#PARTIE 2

#MEME MODELE MAIS SOUS CONTRAINTE

mod_c2 <- lm(Goals~C(Club, base =2), data)
summary(mod_c2)

# on a mis alpha égale a 0, ducoup le Club de l'AC Milan est passer en Intercept.
#Ce qui a compare les estimations des nombres de buts
#des différents clubs en la comparant a celle de l'AC Milan.
#On observe donc que les club ayant une estimation positive on marqué plus de but en moyenne que l'AC Mil
#Tandis que ceux qui ont une estimation négative on marqué moins de but en moyenne que l'AC Milan.
#Mais si l'on veut déduire que la différence de but entre l'AC Milan et un des clubs est significative on
# On voit donc que le Fc Barcelone a marqué en moyenne 2,2 but par match de plus que l'AC Milan et la p-v
#On peut donc affirmer que la différence de but marqué en moyenne entre le Fc Barcelone et l'AC Milan est
# Dans le cas inverse si on s'intéresse au club de Cordoba on peut voir que l'estimation de but moyen mar
#Puis la p-valeur est de 0.03 donc on rejette h0 est on peut affirmer que la différence de but entre cord

#1 la contrainte est alpha=0 intercept nest plus la

#estimation des espérance de Y
summary(mod_c2)$coefficients[,1]

# Tester au risque 5
```

2) Extrait de sorti de code réalisé :

Fiche bilan SAE

```
##
## Call:
## lm(formula = Goals ~ Club - 1, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.093 -1.614 -1.264  0.425 45.907
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## ClubAC Ajaccio      1.2644     0.3627   3.486 0.000492 ***
## ClubAC Milan        1.8916     0.2144   8.822 < 2e-16 ***
## ClubAlaves          1.6400     0.6767   2.424 0.015373 *
## ClubAlmeria         1.4571     0.3302   4.413 1.02e-05 ***
## ClubAngers          1.3621     0.4443   3.066 0.002172 **
## ClubArles-Avignon    0.6563     0.5981   1.097 0.272548

```

m0 <- lm(Goals~0, data = data)
anova(m0, mod_foot)

	Res.Df <dbl>	RSS <dbl>	Df <dbl>	Sum of Sq <dbl>	F <dbl>	Pr(>F) <dbl>
1	21551	318760.0	NA	NA	NA	NA
2	21389	244835.5	162	73924.52	39.86481	0

2 rows