

STATISTICAL ANALYSIS AND DOCUMENT MINING



Pedro L. C. Rodrigues
pedro.rodrigues@inria.fr



Romain Rombourg
romain.rombourg@univ-grenoble-alpes.fr



Daria Bystrova
daria.bystrova@inria.fr

Website: chamilo.grenoble-inp.fr/courses/ENSIMAG4MMSADM

STATISTICAL ANALYSIS AND DOCUMENT MINING

TP

- *R* programming language
- Groups of **three** students
- Four or five TPs with reports to be uploaded on TEIDE
- You should write them in **english** and using **.rmd**

Final grade

Half the average of scores on each TP report

Half the score at final exam

Presence in CM, TD, TP is required and evaluated

Call:

```
lm(formula = sales ~ TV + radio + newspaper, data = adv)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-8.8277	-0.8908	0.2418	1.1893	2.8292

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	2.938889	0.311908	9.422	<2e-16	***
TV	0.045765	0.001395	32.809	<2e-16	***
radio	0.188530	0.008611	21.893	<2e-16	***
newspaper	-0.001037	0.005871	-0.177	0.86	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.686 on 196 degrees of freedom

Multiple R-squared: 0.8972, Adjusted R-squared: 0.8956

F-statistic: 570.3 on 3 and 196 DF, p-value: < 2.2e-16

Call:

```
lm(formula = sales ~ TV + radio + newspaper, data = adv)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-8.8277	-0.8908	0.2418	1.1893	2.8292

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	2.938889	0.311908	9.422	<2e-16	***
TV	0.045765	0.001395	32.809	<2e-16	***
radio	0.188530	0.008611	21.893	<2e-16	***
newspaper	-0.001037	0.005871	-0.177	0.86	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.686 on 196 degrees of freedom

Multiple R-squared: 0.8972, Adjusted R-squared: 0.8956

F-statistic: 570.3 on 3 and 196 DF, p-value: < 2.2e-16

Call:

```
lm(formula = sales ~ TV + radio + newspaper, data = adv)
```

Residuals:

Min	1Q	Median	3Q	Max
-8.8277	-0.8908	0.2418	1.1893	2.8292

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.938889	0.311908	9.422	<2e-16 ***
TV	0.045765	0.001395	32.809	<2e-16 ***
radio	0.188530	0.008611	21.893	<2e-16 ***
newspaper	-0.001037	0.005871	-0.177	0.86

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.686 on 196 degrees of freedom

Multiple R-squared: 0.8972, Adjusted R-squared: 0.8956

F-statistic: 570.3 on 3 and 196 DF, p-value: < 2.2e-16

Call:

```
lm(formula = sales ~ newspaper, data = adv)
```

Residuals:

Min	1Q	Median	3Q	Max
-11.2272	-3.3873	-0.8392	3.5059	12.7751

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	12.35141	0.62142	19.88	< 2e-16 ***
newspaper	0.05469	0.01658	3.30	0.00115 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.092 on 198 degrees of freedom

Multiple R-squared: 0.05212, Adjusted R-squared: 0.04733

F-statistic: 10.89 on 1 and 198 DF, p-value: 0.001148

	TV	radio	newspaper	sales
TV	1.000000000	0.05480866	0.05664787	0.7822244
radio	0.05480866	1.000000000	0.35410375	0.5762226
newspaper	0.05664787	0.35410375	1.000000000	0.2282990
sales	0.78222442	0.57622257	0.22829903	1.00000000

Call:

```
lm(formula = life ~ rpm + brand, data = df)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.5527	-1.7868	-0.0016	1.8395	4.9838

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	36.98560	3.51038	10.536	7.16e-09	***
rpm	-0.02661	0.00452	-5.887	1.79e-05	***
brandB	15.00425	1.35967	11.035	3.59e-09	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.039 on 17 degrees of freedom

Multiple R-squared: 0.9003, Adjusted R-squared: 0.8886

F-statistic: 76.75 on 2 and 17 DF, p-value: 3.086e-09

Call:

```
lm(formula = life ~ rpm + brand, data = df)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.5527	-1.7868	-0.0016	1.8395	4.9838

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	36.98560	3.51038	10.536	7.16e-09	***
rpm	-0.02661	0.00452	-5.887	1.79e-05	***
brandB	15.00425	1.35967	11.035	3.59e-09	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.039 on 17 degrees of freedom

Multiple R-squared: 0.9003, Adjusted R-squared: 0.8886

F-statistic: 76.75 on 2 and 17 DF, p-value: 3.086e-09

```
> contrasts(df$brand)
      2
A 0
B 1
```