

Cours de Détection de Variables et tests multiples

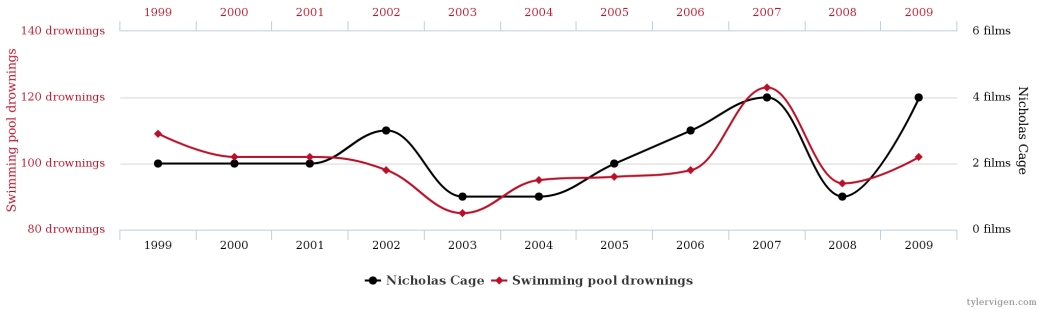
Patricia Reynaud-Bouret

Année 2016-2017



Mise en jambe

Number of people who drowned by falling into a pool
correlates with
Films Nicolas Cage appeared in



Qu'en pensez-vous ? (cf site <http://tylervigen.com/spurious-correlations>)

I Correlation

$$r = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad \text{ou} \quad \text{Cov}(X, Y) = E(XY) - E(X)E(Y)$$

$$\sigma_X = \sqrt{\text{Var} X}, \quad \sigma_Y = \sqrt{\text{Var} Y}$$

Version empirique

Si on a $(x_1, y_1), \dots, (x_n, y_n)$ iid on peut estimer r par

$$\hat{r} = \frac{(\frac{1}{n} \sum_i x_i y_i) - \bar{x} \bar{y}}{\sqrt{\frac{1}{n} \sum_i x_i^2 - \bar{x}^2} \sqrt{\frac{1}{n} \sum_i y_i^2 - \bar{y}^2}} \quad \text{NB: } \hat{r} \in [-1, 1]$$

La "agence populaire" va dire qqch comme si $|\hat{r}| > 0.5$ dépendance forte entre les variables.

RStudio File Edit Code View Plots Session Build Debug Tools Window Help

lin_simple.R

```

1 setwd("~/Desktop/travail/cours/M2lmea")
2
3 stupid=read.table("Data1.dat", header=TRUE)
4
5 cor(stupid$Drawings, stupid$Cage.Mov)
6

```

nom des colonnes du Dataframe stupid

pour se mettre dans le bon répertoire

pour lire les données ne pas oublier l'entête ou pas si pas sûr → readLines

Environment History

Global Environment

Data

- stupid 11 obs. of 3 variables
- stupid.data 11 obs. of 3 variables

Files Plots Packages Help Viewer

R: Correlation, Variance and Covariance (Matrices)

cor (stats)

R Documentation

Correlation, Variance and Covariance (Matrices)

Description

var, cov and cor compute the variance of x and the covariance or correlation of x and y if these are vectors. If x and y are matrices then the covariances (or correlations) between the columns of x and the columns of y are computed.

cov2cor scales a covariance matrix into the corresponding correlation matrix *efficiently*.

Usage

```

var(x, y = NULL, na.rm = FALSE, use)
cov(x, y = NULL, use = "everything",
    method = c("pearson", "kendall", "spearman"))
cor(x, y = NULL, use = "everything",

```

Console

```

~/Desktop/travail/cours/M2lmea
5 2003 1 85
6 2004 1 95
7 2005 2 96
8 2006 3 98
9 2007 4 123
10 2008 1 94
11 2009 4 102
> help(cor)
No documentation for 'cor' in specified packages and libraries:
you could try '??cor'
> ??correlation
> cor(stupid$Drawings, stupid$Cage.Mov)
[1] 0.6660043

```

quand on a oublié, on demande !!

Une corrélation de 0.66 → fort ça doit bien dépendre l'un de l'autre!

A ce stade, vous devez vous dire "on est en train de faire l'impaté qui, il ne peut pas y avoir de liens entre les hayrides et les films de Cage."

NB spurious = fallacieux

Quelles explications voyez-vous?? Réponses possible:

- Les bases de la sagesse populaire ne veulent rien dire
- C'est qu'il y a une bonne base, que veut-on en faire?
- Il (Tyler Vigen) l'a cherché! / a triché??

↳ on va déconstruire les 3 points!!

II Corrélation, indépendance et lien de cause à effet

1/ Des trucs manifestement dépendants et qui ne se corrélaient pas


2/ Des trucs manifestement complètement "indépendants" et qui se corrélaient

Chercher Comment simuler qqch comme ça ?

Pon 1 / Si on est indépendant (au vrai sens math)

Alors $E(XY) = E(X)E(Y)$ et donc $r = 0$, si bcp de données $\hat{r} \approx 0$

Mais on peut très bien imaginer des trucs dépendants qui peuvent vérifier $r = 0$ (juste premier moment)

Ex  + bruit éventuel \rightarrow on vérifie par simulation.

Pon 2 / Il suffit d'être lié à une 3^e variable (par année)

cf toutes les questions corrélation \nearrow en 1^o de l'année

Ex diplôme en sciences / suicides par pendaison!

NB: A ce moment là, \hat{r} pour estimer est en fait assez bnf car pas iid Conditionnellement à l'année.

Cas 1

```

1 setwd("~/Desktop/travail/cours/M2Imea")
2
3 stupid=read.table("Data1.dat", header=TRUE)
4
5 cor(stupid$Drawings,stupid$Cage.Mov)
6
7 ##### Verif par simu de variables dependantes non correlees
8
9
10 n= 20
11 sigma=0.01
12 x=seq(-1,1,length.out=n)
13 y=x^2+rnorm(n,sd=sigma)
14
15 plot(x,y)
16
17 cor(x,y)
  
```

← années - nous à faire varier sigma la corrélation va augmenter avec le bruit!!

```

> plot(x,y)
> cor(x,y)
[1] -1.704156e-16
> n= 20
> sigma=0.01
> x=seq(-1,1,length.out=n)
> y=x^2+rnorm(n,sd=sigma)
> plot(x,y)
> cor(x,y)
[1] -0.009980623
  
```

← qd $\sigma = 0!!$

Corrélation ridicule!!

Environment History

Global Environment

Data

- stupid 11 obs. of 3 variables
- stupid.data 11 obs. of 3 variables

Values

```

n      20
sigma  0.01
x      num [1:20] -1 -0.895 -0.789 -0.684 -0.579 ...
y      num [1:20] 1.014 0.828 0.62 0.451 0.337 ...
  
```

Files Plots Packages Help Viewer

Le graphe montre clairement la dépendance!!

Cas 2

```

19 year = 1999:2009
20 n=length(year)
21
22 sigma=0.1
23 a1= 0.3
24 a2=100
25
26
27 bruit1= rnorm(n,sd=sigma)
28 bruit2=rnorm(n,sd=sigma)
29 bruit1=bruit2
30
31 y1= 0.3 * year + bruit1
32 y2=100* year + bruit2
33
34 plot(year,y1,col="red",axes=FALSE,ylab="")
35 axis(side=4,col="red",ylab="y1")
36 par(new=TRUE)
37 plot(year,y2)
38
39
40 cor(y1,y2)
  
```

← tirages aléatoires!
← ce n'est pas les mêmes!

```

> bruit1=bruit2
[1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
> y1= 0.3 * year + bruit1
> y2=100* year + bruit2
> plot(year,y1,col="red",axes=FALSE,ylab="")
> axis(side=4,col="red",ylab="y1")
> par(new=TRUE)
> plot(year,y2)
> cor(y1,y2)
[1] 0.9967143
  
```

Environment History

Global Environment

```

bruit1 num [1:11] -0.096 0.0756 0.0825 0.1285 0.0379 ...
bruit2 num [1:11] 0.0258 -0.1175 0.014 0.0153 0.025 ...
n      11L
sigma  0.1
x      num [1:20] -1 -0.895 -0.789 -0.684 -0.579 ...
y      num [1:20] 1.014 0.828 0.62 0.451 0.337 ...
y1     num [1:11] 600 600 600 601 601 ...
y2     num [1:11] 2e+05 2e+05 2e+05 2e+05 2e+05 ...
year   int [1:11] 1999 2000 2001 2002 2003 2004 2005 2006 ...
  
```

Files Plots Packages Help Viewer

On voit à une dépendance mais clairement il conditionnellement à l'année.

Corrélation énorme!!

3/ L'indépendance

Sans rentrer dans tous les détails précises vous devez en connaître au moins un ??

↳ Chi-deux d'indépendance

(X_i, Y_i) ind $i=1..n$ et catégoriel valeur de $1..r$ pour X
 $1..s$ pour Y .

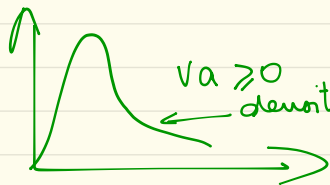
Alors si les $X_i \perp\!\!\!\perp Y_i$,

N_{jk} nb de couples qui valent (j, k) (ie $X_i=j, Y_i=k$)
 $N_{j.}$ nb de couples tq $X_i=j$
 $N_{.k}$ nb de couples tq $Y_i=k$

les f^o de rép cv

Alors
$$\sum_{j,k} \frac{(N_{jk} - \frac{N_{j.} \cdot N_{.k}}{n})^2}{\frac{N_{j.} \cdot N_{.k}}{n}} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi^2_{(r-1)(s-1)}$$

nb observé (pointing to N_{jk})
nb attendu si indep... (pointing to $\frac{N_{j.} \cdot N_{.k}}{n}$)



bi du χ^2 à (ici) $(r-1) \times (s-1)$ degré de liberté

Pour que approx valable il faut que tous les nb attendus ≥ 5 .

Comme c'est pour des va catégorielles comment fait on ici?

→ Cage mor OK (a priori)

→ Drownings (on pourrait mais on va avoir du mal pour $N_{at} \geq 5 \rightarrow$ regrouper)

`data=stupid[, -1]` ← j'enlève la colonne année

`tab1=table(data)` #matrice de contingence

`chi1=chisq.test(tab1)` # pval : 15 % (donc Indep acceptee) mais regardez le warning

`chi1$expected` # pas bon, pas ≥ 5 (ça explique le warning)

`data2=as.data.frame(cbind(data[,1],(data[,2]>100)))` # j'ai fait deux catégories sur les drownings
`tab2=table(data2)`

`chi2=chisq.test(tab2)`

`chi2$expected` # toujours pas bon

`data3 =as.data.frame(cbind(data[,1]>2,(data[,2]>100)))` # ~~2~~ deux catégories sur les deux valeurs

`tab3=table(data3)`

`chi3=chisq.test(tab3)`

`chi3$expected` # toujours pas bon (allez voir ce que c'est que la Yates correction dont R parle)

en fait quand on arrive a des tables 2×2 on peut faire le test exact de fisher (Allez voir ce que c'est si vous ne savez pas....)

`fisher.test(tab3)` # et la on ne rejette rien du tout, pval 1 ...

⚠ A force de mettre des catégories de plus en plus grosses on perd de l'information

⚠ Un test est fait pour accepter son hypothèse nulle (ici \perp entre X et Y), il ne va donner des petites p-values qu'à contre cœur... donc pas d'info si pval grande!!

Pour les données parabol → pas catégorielles → on fait les catégories nous même

```
n= 100 # j'augmente un peu pour pas être gêné par les >=5
```

```
sigma=0.01
```

```
x=seq(-1,1,length.out=n)
```

```
y=x^2+rnorm(n,sd=sigma)
```

```
plot(x,y)
```

```
cor(x,y) # corrélation ridicule, les variables sont pas corrélées.
```

```
data.para=as.data.frame(cbind((x>-0.5)+(x>0.5),(y>0.25)))
```

```
tab.para=table(data.para)
```

```
chi.para=chisq.test(tab.para) # pas de message d'erreur
```

```
chi.para # p-valeur clairement ridicule, on rejette l'indépendance, les variables sont dépendantes
```

```
chi.para$expected # tout va bien pour l'approximation
```

Essayez tout seul sur les var qui croissent linéairement en % de l'année → ça dira dépendance

Car y_1 et y_2 sont \perp conditionnellement à l'année mais pas indépendamment.

⇒ Le problème de la détection de variables est que c'est possible si on observe bien toutes les variables mais si des variables sont cachées on va trouver des liens absurdes → toujours faire attention avec les analyses.

4/ Cause à effet

Rien n'est plus trompeur que de confondre corrélation et lien de cause à effet.

cause à effet : - au moins lien temporel : cause avant effet (et pas juste simultanée)
- au mieux une fois une dépendance soupçonnée on si on peut faire l'expérience de manipuler qu'une variable et voir l'effet

III et la modélisation dans tout ça ?

Sur les stupid data, c'est pas parce que le test d'II accepte que ça dit quoique a tort (un test qui accepte ne dit rien en fait)

↳ Y aurait-il un modèle statistique dans lequel les bones demeurent sur \hat{f} prendrait un sens ?

1/ le modèle de régression linéaire

$$Y_i = a_0 + b_0 X_i + \varepsilon_i$$

avec ε_i bruit (ie $E(\varepsilon_i | X_i) = 0$)

a_0 et b_0 inconnus.

Pour deviner a_0 et b_0 , on peut minimiser la distance l^2

↳ on cherche a et b qui minimise $\sum_i (Y_i - (a + bX_i))^2$

↳ Comment fait-on ??

$$g(a, b) = \sum_i (Y_i - (a + bX_i))^2$$

$$= \sum_i (Y_i^2 - 2(a + bX_i)Y_i + (a + bX_i)^2)$$

$$= \sum_i Y_i^2 - 2a \sum_i Y_i - 2b \sum_i X_i Y_i + a^2 n + 2ab \sum_i X_i + b^2 \sum_i X_i^2$$

$$\frac{\partial}{\partial a} = -2 \sum_i Y_i + 2an + 2b \sum_i X_i \Rightarrow a = \overline{Y} - b \overline{X}$$

min en

$$Q(\hat{a} - b\bar{X}, b)$$

$$= \sum y_i^2 - 2a \sum y_i - 2b \sum x_i y_i + a^2 n + 2ab \sum x_i + b^2 \sum x_i^2$$

$$= \sum y_i^2 - 2n(\bar{y} - b\bar{X})\bar{y} - 2b \bar{X} y n + (\bar{y} - b\bar{X})^2 n + 2(\bar{y} - b\bar{X})bn\bar{X} + b^2 n\bar{X}^2$$

$$= \sum y_i^2 - 2n(\bar{y})^2 + (\bar{y})^2 n + b(2n\bar{X}\bar{y} - 2n\bar{X}\bar{y}) + b^2(n\bar{X}^2 - n\bar{X}^2)$$

$$\frac{\partial}{\partial b} = 2n(\bar{X}\bar{y} - \bar{X}\bar{y}) + 2bn(\bar{X}^2 - \bar{X}^2)$$

\Rightarrow min (vérifier monotonie / chgt de signe dans le bon sens)

$$\hat{b} = \frac{\overline{XY} - \bar{X}\bar{Y}}{\overline{X^2} - \bar{X}^2} = \frac{\text{Cov}_n(X, Y)}{\text{Var}_n(X)}$$

$$\hat{a} = \bar{Y} - \hat{b}\bar{X}$$

Lien avec la Corrélation

$$\frac{\text{Var expliquée par le modèle}}{\text{Var intrinsèque}} = \frac{\sum (\bar{Y} - \hat{y}_i)^2}{\sum (y_i - \bar{Y})^2}$$

$$= \frac{\sum (\hat{b}\bar{X} - \hat{b}x_i)^2}{\sum (y_i - \bar{Y})^2} = \frac{\text{Var}_n X \hat{b}^2}{\text{Var}_n Y} = \frac{\text{Cov}_n(X, Y)^2}{\text{Var}_n X \text{Var}_n Y} = \left(\hat{r}\right)^2 \quad (\text{Coeff de Corr. empirique}^2)$$

où $\hat{y}_i =$ les préd du modèle $= \hat{a} + \hat{b}x_i$

Donc \hat{r} est proportionnel à \hat{b} .

⇒ faire un test de la nullité de \hat{b} \Leftrightarrow tester nullité de \hat{r}
et donc la corrélation
éventuelle entre X et Y

2/ Le modèle linéaire gaussien (cas de la régression simple)

Comme on l'a vu avant, tester demande de connaître la loi
(asymptotique éventuellement) de la stat. de test (ici \hat{b})

↳ dépend forcément du bruit sous-jacent ε_i .

Si on suppose ε_i gaussien iid $\mathcal{N}(0, \sigma^2)$ σ^2 éventuellement inconnu

On peut alors connaître la loi de \hat{b} / $b \rightarrow$ test de $b=0$

NB: La plus grosse partie théorique des modèles
linéaires gaussiens consiste à savoir estimer les param.
(ici b) et à donner leur loi !!

(c'est beaucoup de proba. et d'algèbre linéaire
mais dans un seul but: arriver à interpréter ce qui se
passe sur les vraies données.)

Pour le moment, on va faire comme si vous saviez tout
pour voir comment ça s'applique

RStudio File Edit Code View Plots Session Build Debug Tools Window Help

modlin_readdata.R

```

90 ## modele lineaire de regression gaussienne
91
92 par(mfrow=c(1,2)) ← 2 dessins juxtaposés
93 plot(stupid$Drownings, stupid$Cage.Mov)
94 plot(stupid$Cage.Mov, stupid$Drownings)
95
96 res=lm(stupid$Drownings~stupid$Cage.Mov)
97
98 lines(stupid$Cage.Mov, res$fitted.values, type="l", col="red")
99
100 summary(res) ← donne toutes les infos calculées
101
102
103 (Top Level)

```

Environment History

Global Environment

- ccni1 List of 9
- chi.para List of 9
- chi1 List of 9
- chi2 List of 9
- chi3 List of 9
- n 100
- res List of 11
- sigma 0.01
- tab.para 'table' int [1:3, 1:2] 0.49 0.25 1.25

Files Plots Packages Help Viewer

Zoom Export

Call:

```
lm(formula = stupid$Drownings ~ stupid$Cage.Mov)
```

Residuals:

Min	1Q	Median	3Q	Max
-8.418	-6.597	1.045	3.224	12.582

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	87.134	5.443	16.009	6.4e-08 ***
stupid\$Cage.Mov	5.821	2.173	2.678	0.02531 **

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

Residual standard error: 7.585 on 9 degrees of freedom
Multiple R-squared: 0.4436, Adjusted R-squared: 0.3817
F-statistic: 7.174 on 1 and 9 DF, p-value: 0.02527

Handwritten notes in red and purple:

- 2 dessins juxtaposés
- commande pour faire calcul mod. lin.
- $y = a + b \cdot x$ avec a et b inconnus.
- valeurs prédites \hat{y}_i
- donne toutes les infos calculées
- suppose bruit gaussien
- ici plus légitime
- préal de $b=0$

suppose bruit gaussien
dans $Cage.mov = a + b \cdot Drownings + \text{bruit}$
 \Rightarrow idiot $Cage.mov$ prend 4 valeurs!
mais $Drownings = a + b \cdot Cage.mov + \text{bruit}$
est plus raisonnable

Retien Code des pval (***) etc)

ici le test de $b=0$ fait automatiquement pval: 0.02

↳ test accepte au niveau 1%
rejeté au niveau 5% → le lien semble avéré??

Oui mais

1/ et si pas Gaussien??

test de Shapiro et Wilk (on refra...)

(test le + puissant pour vérifier normalité!)

Shapiro test (residuals) \rightarrow pval 26% semble Gaussien

! Mumble... (bien sûr qu'un test de Kolmogorov pas Gaussien est disant, mais approx semble de et approx Binomiale / Gaussien)

2/ Donc y a un lien!!??!!

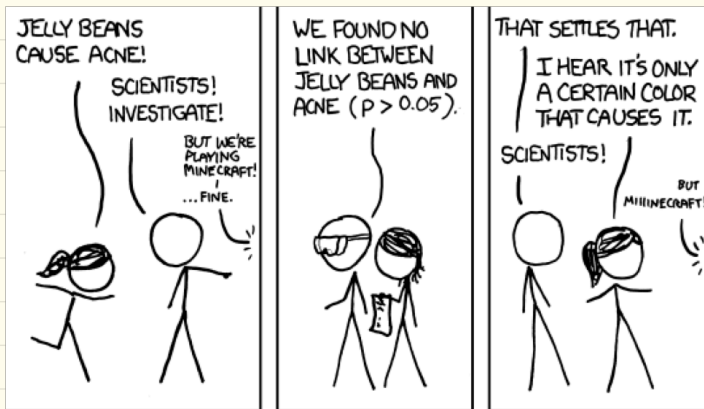
! A l'interprétation d'une p-value en particulier en cas de tests multiples non dits!

Prop pval

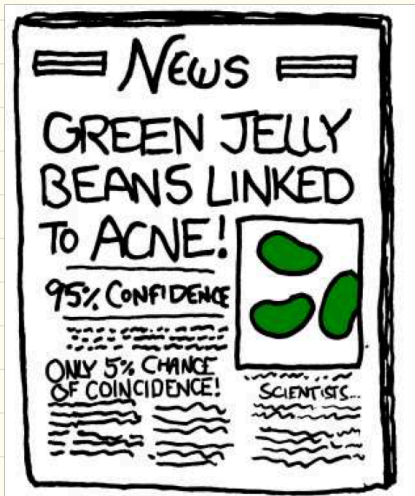
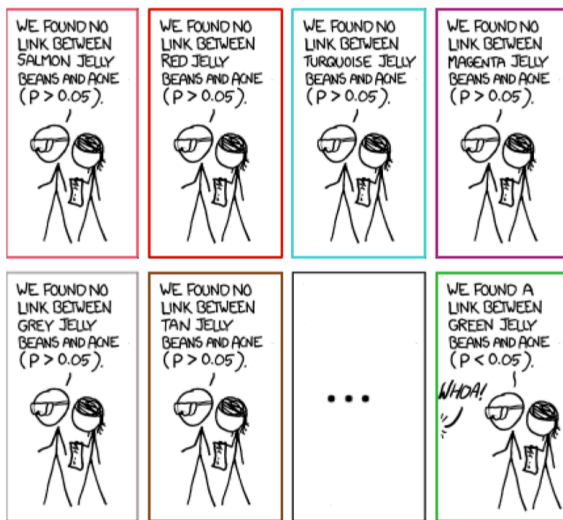
Smooth, $P(\text{pval} \leq \alpha) \approx \alpha$ (on y reviendra)

Si $\alpha = 5\%$ (test au niveau 5%)

$\text{pval} \leq 5\%$ a 5% de chances de se produire!!



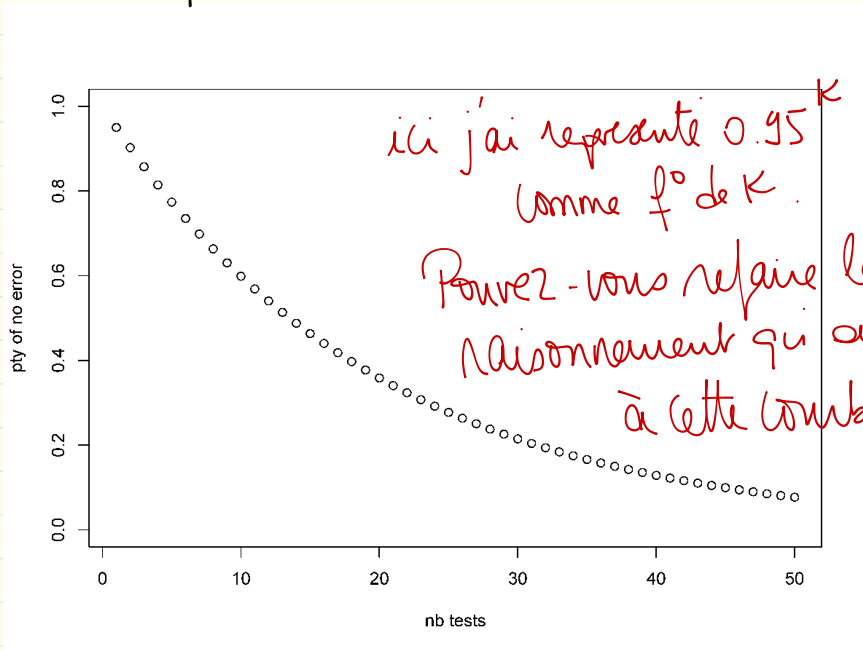
source : xkcd



Or qu'a fait Tyler Vigen ? Combien a-t-il testé de courbes bizarres avant de trouver la paire qui va avec les films de Cage ???

Si chaque donnée est vraiment indépendante des films de Cage et qu'on teste K données différentes, on trouve que la probabilité de trouver un lien alors qu'il n'y en a pas est

↑
ne pas



Que peut-on dire ?

Bonferroni (et on revient pourquoi) dit qu'il ne faut pas comparer à α mais α/K

avec $p_{val} = 2\%$ il suffit de 3 tests et $\alpha = 5\%$

pour que ça soit plus significatif !!
→ Bref, "Vigen a triché..." en ne disant pas combien de données il a testé !!

Tests non paramétriques

I Représentations graphiques et jugement qualitatif

1. Fonction de répartition

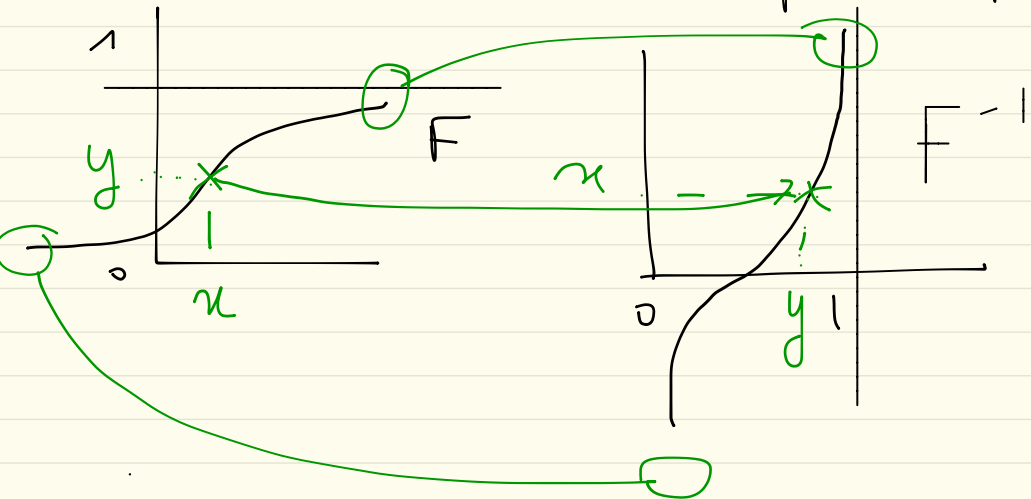
- la Fonction de répartition $t \mapsto P(X \leq t) = F(t)$ caractérise la loi de la variable réelle X .

- F est une fonction \nearrow $\lim_{-\infty} F = 0$; $\lim_{+\infty} F = 1$

- La fonction quantile $F^{-1}(t) = \inf\{x / F(x) \geq t\}$

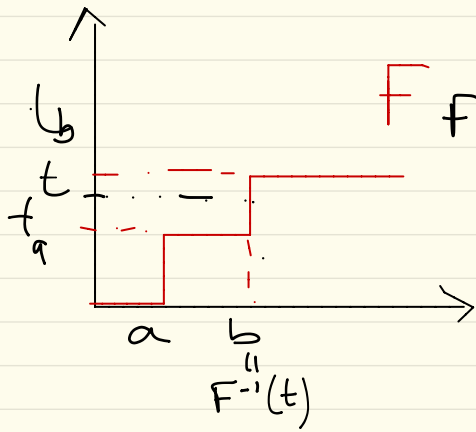
• Si X a densité $f > 0$, F est bijective et continue

F^{-1} est la fonction réciproque

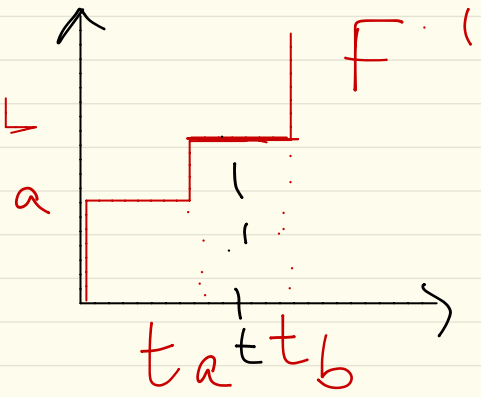


Si non F^{-1} est quand même bien défini

NB: Si $U \sim \mathcal{U}[0,1]$, $F^{-1}(U) \sim X$
 & $F \circ F^{-1} = \text{id}$, $F(X) \sim \mathcal{U}[0,1]$



$$F \circ F^{-1}(t) = t$$



• Si X_1, \dots, X_n iid de c. d. f F
 alors la fonction de répartition empirique est donnée par

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq t}$$

C'est un estimateur sans biais et consistant de F

- Sans biais : $\forall t \in \mathbb{R} \quad E[\hat{F}_n(t)] = F(t)$

- Consistant (dans l'espace des fonctions $\mathcal{F} = \{F \uparrow \lim_{t \rightarrow -\infty} = 0, \lim_{t \rightarrow +\infty} = 1\}$ muni de la $\|\cdot\|_b$)

$\sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \xrightarrow[n \rightarrow +\infty]{\text{P.S.}} 0$ (Théorème de Glivenko-Cantelli)

2. Représentation graphique et Qq plot

a) Si on pense que les X_i suivent une certaine loi de c.d.f. F_0

on peut toujours tracer \hat{F}_n et F_0 et voir si ils sont proches

⇒ Utiliser la commande `ecdf` de R
Simuler n va gaussiennes i.i.d $CP(0,1)$
Vérifier le fit avec F_0 cdf de la $CP(0,1)$
regarder comment ça se déforme si on ne met pas la bonne moyenne ou la bonne variance

⇒ Si on veut vérifier que ça marche mais on ne connaît pas moyenne et variance

↳ Que suggérez-vous?

Estimer moyenne, variance, tracer F_0 la cdf d'une $CP(\hat{m}, \hat{\sigma}^2)$

b) Version estimation de densité ça vous dit quoi?
↳ histogramme
↳ estimateur à noyau

Que savez-vous sur les histogrammes?

- en pratique : la 1^o hist, la 1^o densité
- choix du pas ?

`n=20 # 100`
`X=rnorm(n)`

avec 100 vous venez qqch de convaincant
mais si n nombre d'observation peut moins convaincant

fonction de repartition

`plot.ecdf(X)`

`abs=seq(-3,3,0.1)`

`lines(abs,pnorm(abs),col='red')`

`lines(abs,pnorm(abs,mean=0.5),col='green')`

`lines(abs,pnorm(abs,sd=0.5),col='blue')`

`lines(abs,pnorm(abs,mean=mean(X),sd=sd(X),col='orange')`

`legend('topleft',c('empirical`

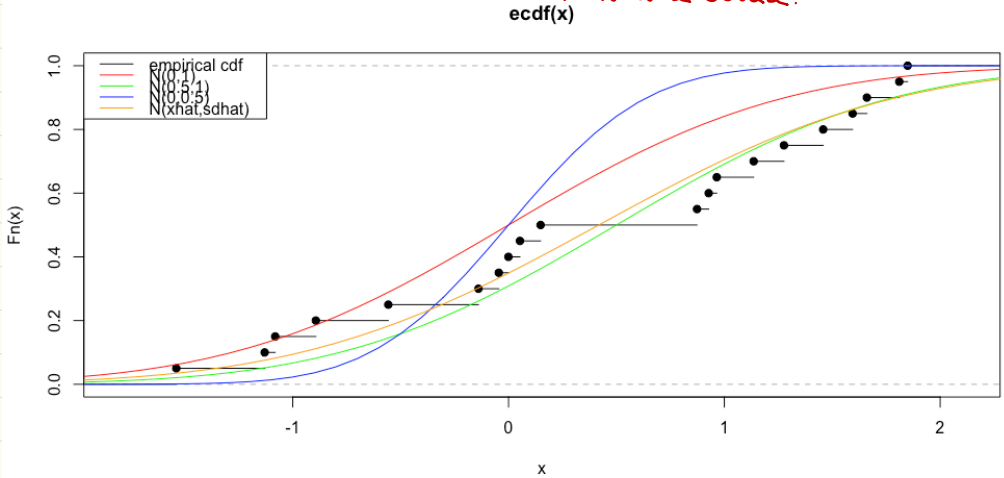
`cdf', 'N(0,1)', 'N(0.5,1)', 'N(0,0.5)', 'N(xhat,sdhat)'), col=c('black', 'red', 'green', 'blue', 'orange'), lty=c(1,1,1,1,1))`

regarder l'aide pour voir ce qu'est ecdf

regardez bien l'aide pour mettre au bon endroit.

sur une seule ligne de commande

La courbe ci dessous a été tracée avec $n=20$
On voit bien que la proximité avec la bonne courbe (rouge)
douteux, bcp + proche de l'estimé (qu'en pensez vous?)
et d'une mauvaise courbe.



estimation de densité

hist(X,freq=FALSE)

rug(X)

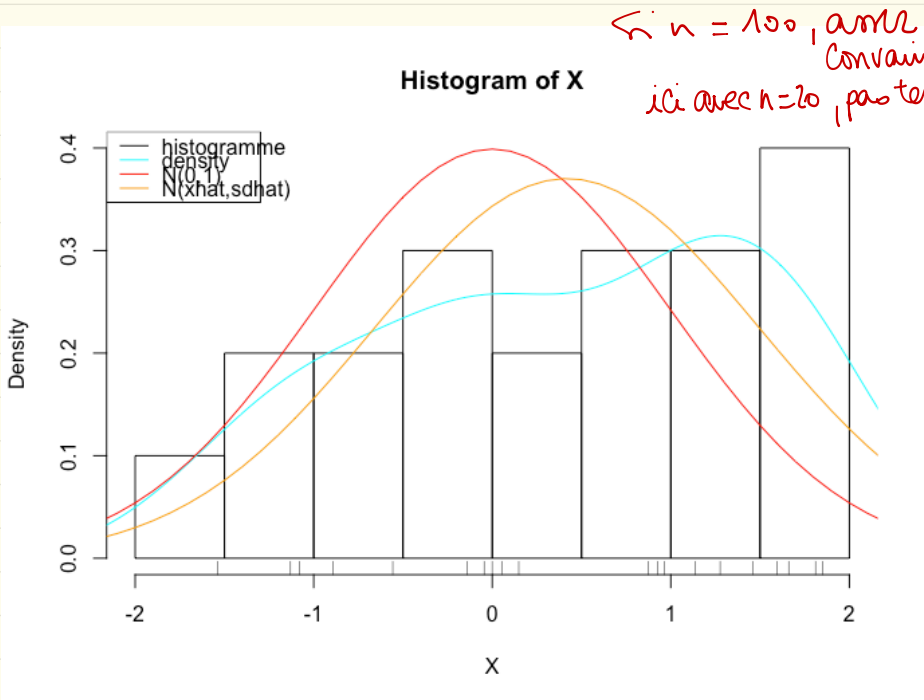
lines(density(X),col='cyan')

lines(abs,dnorm(abs),col='red')

lines(abs,dnorm(abs,mean=mean(X),sd=sd(X)),col='orange')

legend('topleft',c('histogramme','density','N(0,1)','N(xhat,sdhat)'),col=c('black','cyan','red','orange'),lty=c(1,1,1,1))

regarder bien l'aide, voyez règle de Sturges
pour voir où sont les points
attention ce n'est pas la vraie densité mais
un estimateur, dissez l'aide pour voir
ce que c'est!!



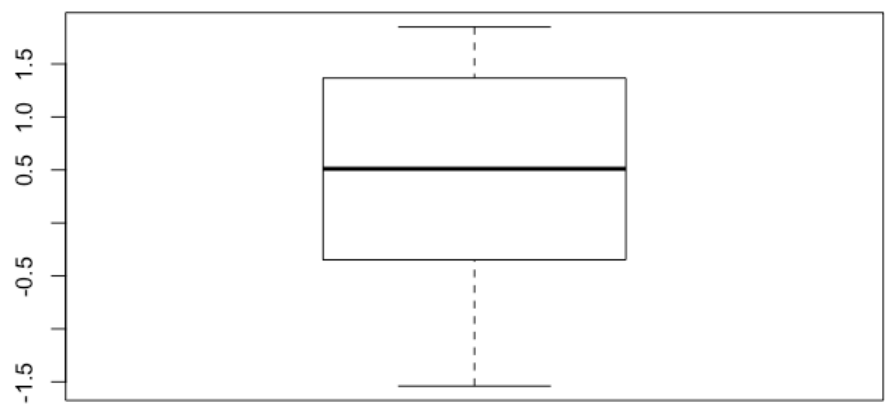
n = 100, assez convaincant
ici avec n=20, pas terrible

c) ### boxplot (boite a moustache)

boxplot(X)

Qd n'ent, finalement boxplot assez suffisant

- Vérifier que vous savez ce qui est tracé
- Regarder ce qui se passe quand on augmente n et en particulier l'apparition d'outliers



d) Qq plot

Idee du qq plot : on est visuellement plus sensible à "voir une droite"

→ Qq plot entre un échantillon X_1, \dots, X_n et une loi (cdf F_0)

$(x, y) \in \text{Qq plot}$ si $x = F_0^{-1}(t)$ et $y = \hat{F}_n^{-1}(t)$
↑ vrai quantile ↑ quantile empirique

Si $X_1, \dots, X_n \sim F_0$ le qq plot devrait s'aligner sur la diagonale $x=y$

↳ qq plot entre deux échantillons

$x_1 \dots x_n$ cdf empirique \hat{F}_n , $y_1 \dots y_n$ cdf empirique \hat{G}_n

$(x, y) \in \text{qq plot}$ si $\exists t$, $x = \hat{F}_n^{-1}(t)$ et $y = \hat{G}_n^{-1}(t)$

Si les deux échantillons ont même loi ça devrait s'aligner sur la diagonale $x=y$

↳ qq plot pour loi normale

Si $X \sim \mathcal{N}(m, \sigma^2)$ alors $\frac{X-m}{\sigma} \sim \mathcal{N}(0, 1)$

donc si Φ cdf de la $\mathcal{N}(0, 1)$

$$P(X \leq t) = P\left(\frac{X-m}{\sigma} \leq \frac{t-m}{\sigma}\right) = \Phi\left(\frac{t-m}{\sigma}\right)$$

donc si $F^{-1}(u)$ quantile de X

$$F(F^{-1}(u)) = u = \Phi\left(\frac{F^{-1}(u)-m}{\sigma}\right) \text{ et } F^{-1}(u) = m + \sigma \Phi^{-1}(u)$$

Donc si $X \sim \mathcal{N}(m, \sigma^2)$, les quantiles empiriques sont approximativement une fonction linéaire de $\Phi^{-1}(u)$
(quantile gaussien centré réduit)

qq plot pour la loi normale

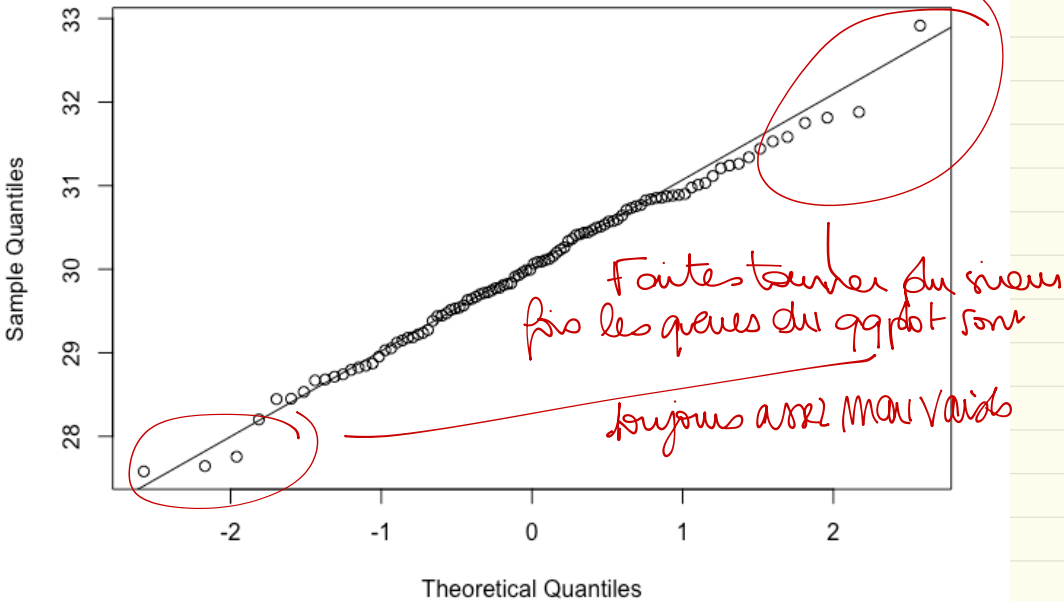
```
X=rnorm(100,mean=30)
```

```
qqnorm(X,main='qqplot pour loi normale')
```

```
qqline(X)
```

← regarder l'aide ça rajoute la ligne qui passe par les 2 pts 1^{er} et 3^e quantiles (ie $t=0.25$ et $t=0.75$)

qqplot pour loi normale



```
### qq plot de deux data sets
```

```
X=rnorm(100)
```

```
nb=rbinom(200,0.25)
```

```
Y=c(rnorm(nb,mean=3,sd=0.4),rnorm(200-nb))
```

```
Z=rnorm(100)
```

```
qqplot(X,Y,xlim=c(-2,2),ylim=c(-3,5),main='qqplot avec deux datasets')
```

```
par(new=TRUE)
```

```
qqplot(X,Z,xlim=c(-2,2),ylim=c(-3,5), col='red')
```

```
qqline(X,distribution=function(p) qnorm(p,mean=0,sd=1))
```

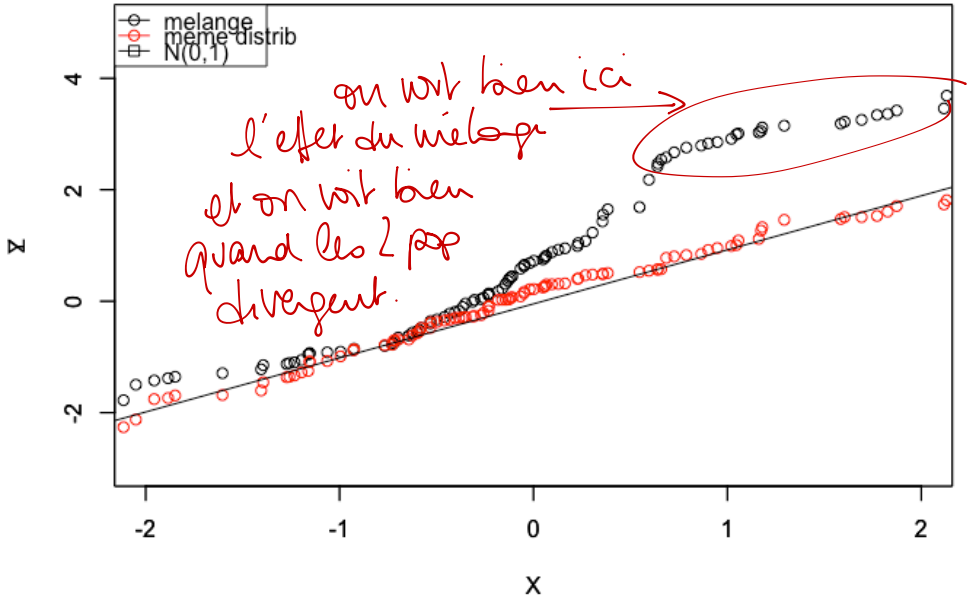
```
legend('topleft',c('melange','meme distrib','N(0,1)'),pch=c(1,1,0),lty=c(1,1,1), col=c('black','red','black'))
```

← quelle loi ai-je simulé ??

← qq plot à 2 échantillon

← je trace la dti qui passe par les 2 quantiles de la v. réelle distib.

qqplot avec deux datasets



II Tests de Kolmogorov Smirnov et variantes

1. Test de KS à un échantillon

On observe X_1, \dots, X_n iid et on se donne une c.d.f. F_0

On veut tester H_0 : les X_i ont pour cdf F_0 contre ce n'est pas le cas
cas particulier de tests plus généraux appelés goodness-of-fit

On va rejeter si la distance entre la cdf empirique et F_0 est trop grande

le test de KS consiste à prendre comme distance (et donc statistique de test)

empirical o.d.f.

$$D_n^{F_0} = \sup_{t \in \mathbb{R}} \left| \hat{F}_n(t) - F_0(t) \right|$$

cdf supposée qu'on veut tester

• Sous H_0 , les X_1, \dots, X_n ont pour cdf F_0 et on peut calculer/simuler la loi de $D_n^{F_0} \Rightarrow$ quantiles $q_n^{F_0}(t)$ et f de rep^o $Q_n^{F_0}(t)$

On rejette si $D_n > q_n(1-\alpha)$ et cela fait un test de niveau α

p valeur : $1 - Q_n^{F_0}(D_n)$ NB : pour simplifier, je fais comme si $Q_n^{F_0} \leftarrow$

• La statistique $D_n^{F_0}$ est calculable facilement

en effet, si on range l'échantillon \rightarrow les statistiques d'ordre

$$X_{(0)} = -\infty \leq X_{(1)} \leq \dots \leq X_{(n)} \leq X_{(n+1)} = +\infty$$

$$D_n^{F_0} = \sup_{i=0, \dots, n} \sup_{t \in [X_{(i)}, X_{(i+1)}[} \left| \hat{F}_n(t) - F_0(t) \right|$$

↑ sur cet intervalle, c'est juste $\frac{i}{n}$

$$= \sup_{i=0, \dots, n} \sup_{t \in [X_{(i)}, X_{(i+1)}[} \left| \frac{i}{n} - F_0(t) \right|$$

mais $F_0 \nearrow$ sur l'intervalle donc le sup est atteint aux extrémités

$$D_n^{F_0} = \max_{i=0, \dots, n} \max \left[\left| \frac{i}{n} - F_0(X_{(i)}) \right| ; \left| \frac{i}{n} - \lim_{t \nearrow X_{(i+1)}} F_0(t) \right| \right] = F_0(X_{(i+1)}) \text{ si } F_0 \text{ c.}$$

x Si F_0 est C^0 , la loi de $D_n^{F_0}$ ne dépend pas de F_0

(pratique pour la tabulation, c'est cette loi qui est utilisée dans R)
→ warnings si R voit qu'il y a des ex de quo ça ne peut pas être continue)

Preuve (stat vient du calcul prodit car les $F(x_i) \sim U[0,1]$)
ou sinon on utilise l'équivalence suivante:

(vraie $\forall F$ même non continue)

$$F^{-1}(x) \leq t \Leftrightarrow x \leq F(t) \quad \text{pour } \forall F \text{ cdf.}$$

$$\sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq t} - F_0(t) \right| \quad (X_i \text{ iid } \sim F_0)$$

$$\sim \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{F_0^{-1}(U_i) \leq t} - F_0(t) \right| \quad (U_i \text{ iid } \sim U[0,1])$$

$$= \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{U_i \leq F_0(t)} - F_0(t) \right|$$

$$= \sup_{x \in F_0(\mathbb{R})} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{U_i \leq x} - x \right|$$

Si $F_0 C^0$, $F_0(\mathbb{R}) = [0,1]$ et c'est $= D_n^{U[0,1]}$

Si non $F_0(\mathbb{R}) \not\subset [0,1]$ et $\leq D_n^{U[0,1]}$

Donc $P(D_n^{F_0} > t) \leq P(D_n^{U[0,1]} > t)$

Les quantiles devraient donc s'appliquer dans pb.

En pratique dans R, ça ne marche pas car il ne calcule pas la bonne stat de test.

x asymptotiquement, TCL $\sqrt{n} (\hat{F}_n(t) - F_0(t)) \xrightarrow[\text{sous } H_0]{\mathcal{L}} \text{CP}(0, F_0(t) \cdot (1 - F_0(t)))$

\Rightarrow on peut mg $\sqrt{n} D_n \xrightarrow[\text{sous } H_0]{\mathcal{L}}$ une distribution connue et tabulée (F_0 est C^0)

NB par/oi

$D_n = \text{en fait } \sqrt{n} D_n$

(ce on normalise par \sqrt{n} la stat de test)

NB: Kolmogorov (33)

$$P(\sqrt{n} \sup_t |\hat{F}_n(t) - F_0(t)| \geq \alpha) \xrightarrow[\text{sous } H_0]{\mathcal{L}} 2 \sum_{k=1}^{\infty} (-1)^{k+1} e^{-2k^2 \alpha^2}$$

Smirnov (42)

$$P(\sqrt{n} \sup_t (\hat{F}_n(t) - F_0(t)) \geq \alpha) \rightarrow e^{-2\alpha^2}$$

On reconnaît une "queue géométrique"

$D_n^+ = \sup_t (\hat{F}_n(t) - F_0(t))$ while pour

$H_0: \forall t F(t) \leq F_0(t)$ contre $\exists t / F(t) > F_0(t)$

car F cdf des X , F_0 cdf de X_0

si $F(t) \leq F_0(t) \forall t$, alors $\sqrt{\frac{F_0^{-1}(t) - F^{-1}(t)}{F_0^{-1}(t)}}$ pour $U \sim \mathcal{U}[0,1]$ $F^{-1}(U), F_0^{-1}(U)$

ie on peut construire des versions de X et X_0 tq $X \geq X_0$
on dit alors que X est stochastiquement plus petit que X_0 .

donc ici si X et X_0 sont des durées de vie

H_0 : " X_0 meurt avant X " H_1 : \exists une durée où X_0 atteint t plus souvent que X .

Massart (90) $\ln P(\sqrt{n}D_n^+ > \alpha) \leq e^{-2\alpha^2}$ et

$$P(\sqrt{n}D_n > \alpha) \leq 2e^{-2\alpha^2} \quad (*)$$

on peut donc en déduire ce que fait le test basé sur D_n l'alternative $H_1: \exists t, F(t) \neq F_0(t)$

Sous H_1 ,

$$\|F - F_0\|_\infty \leq \|F - \hat{F}_n\|_\infty + \|\hat{F}_n - F_0\|_\infty$$

$\xrightarrow{PS} 0$ qd $n \rightarrow \infty$ (à la vitesse \sqrt{n} voir l'étude précédente)
(Glivenko Cantelli)

$$\text{donc } P\left(\|\hat{F}_n - F_0\|_\infty > \frac{\|F - F_0\|_\infty}{2}\right) \xrightarrow{n \rightarrow \infty} 1$$

Or par $*$ par, en posant $\alpha = \frac{\|F - F_0\|_\infty}{2}$
ie $\alpha = \sqrt{\frac{\log(\frac{2}{\alpha})}{2}}$, $P(D_n > \frac{\alpha}{\sqrt{n}}) \leq \alpha$

$$\text{donc } q_n^{F_0}(1-\alpha) \leq \frac{\alpha}{\sqrt{n}}$$

$$\text{donc } P(\text{test rejette}) = P(\|\hat{F}_n - F_0\|_\infty > q_n^{F_0}(1-\alpha))$$

$$\geq P(\|\hat{F}_n - F_0\|_\infty > \frac{\alpha}{\sqrt{n}})$$

$$\geq (\text{pour } n \text{ assez grand}) P(\|\hat{F}_n - F_0\|_\infty > \frac{\|F - F_0\|_\infty}{2})$$

$$\xrightarrow{n \rightarrow \infty} 1$$

Donc sous H_1 , la puissance du test de KS tend vers 1

verification sous H0 du niveau du test de KS

Nsimu=5000

n=30 # taille de l'echantillon

pval_cont=rep(0,Nsimu) # on initialise

for(i in 1:Nsimu)

{

X=rnorm(n) # l'echantillon gaussien F0 continu

a= ks.test(X,'pnorm')

pval_cont[i]=a\$p.value

}

plot.ecdf(pval_cont,main='repartition des pvaleurs')

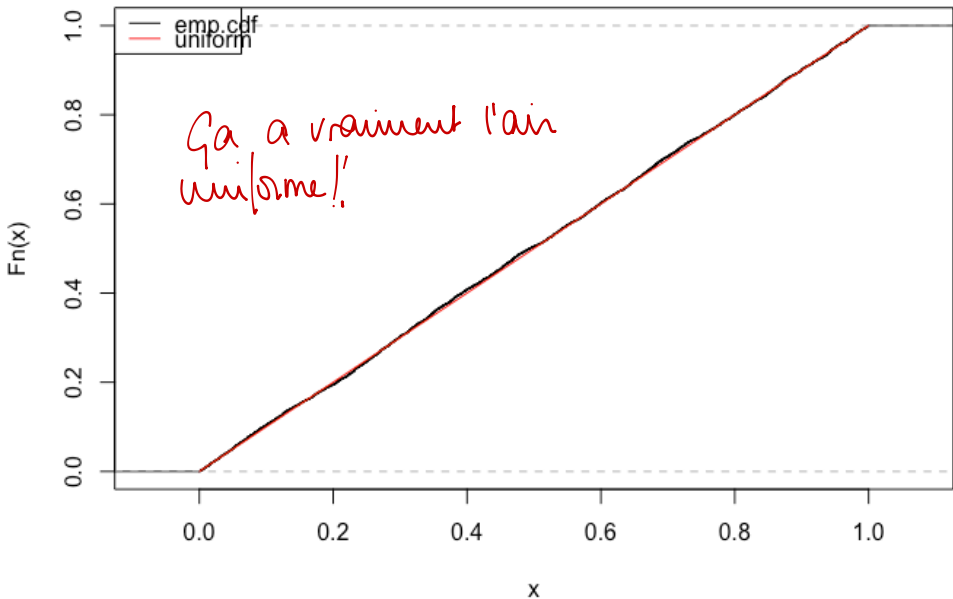
lines(c(0,1),c(0,1),col='red')

legend('topleft',c('emp.cdf','uniform'),col=c('black','red'),lty=c(1,1))

→ Comment le faire? vous?
→ méthode de Monte Carlo
→ jouer avec Nb de simu
→ un lequel on fait KS

NB : amusez vous à prendre
un échantillon d'une distribution
discontinue pour voir
ce qui se passe

repartition des pvaleurs



alpha=0.05 # je verifie que le test est de niveau 5%

niveau_emp=mean(pval_cont<alpha) # comment savoir si la deviation est raisonnable ?

test du "niveau = 5%"

sous H0, niveau_emp est presque une gaussienne $N(0.05, 0.05*0.95/Nsimu)$

erreur_adm= sqrt(0.05*0.95/Nsimu)*qnorm(0.975)

abs(alpha-niveau_emp)>erreur_adm # c'est le rejet de ce test qui est faux donc on accepte

#on peut aussi verifier plus globalement que les pvaleurs sont uniformes

ks.test(pval_cont,'punif')

On teste le test 😊

↳ c'est jamais 5% pile! Regardez la variation en % de Nsimu pourquoi 5000 est plutôt bon??

sous H1

n=100 # taille de l'échantillon

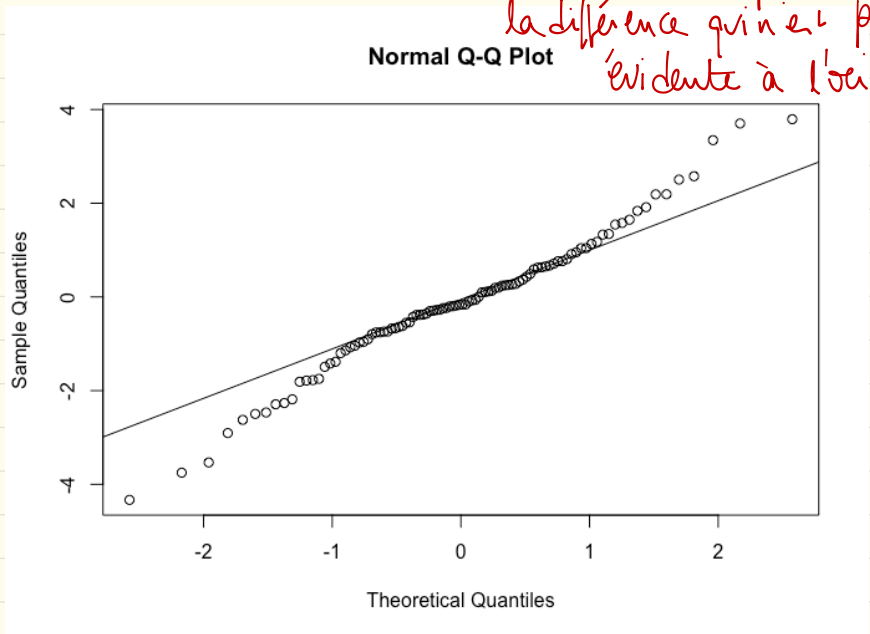
X=(2*rbinom(n,1,0.5)-1)*rexp(n,1) #c'est pas une $N(0,2)$ mais c'est pas loin

qqnorm(X)

qqline(X)

← c'est une loi de Laplace. Derrière sa densité.

↑ test nécessaire pour voir la différence qui n'est pas évidente à l'œil nu



```
hist(X,freq=FALSE)
abs=seq(-6,6,0.1)
lines(abs,dnorm(abs, sd=sqrt(2)),col='red')
```

différence toujours pas flagrante



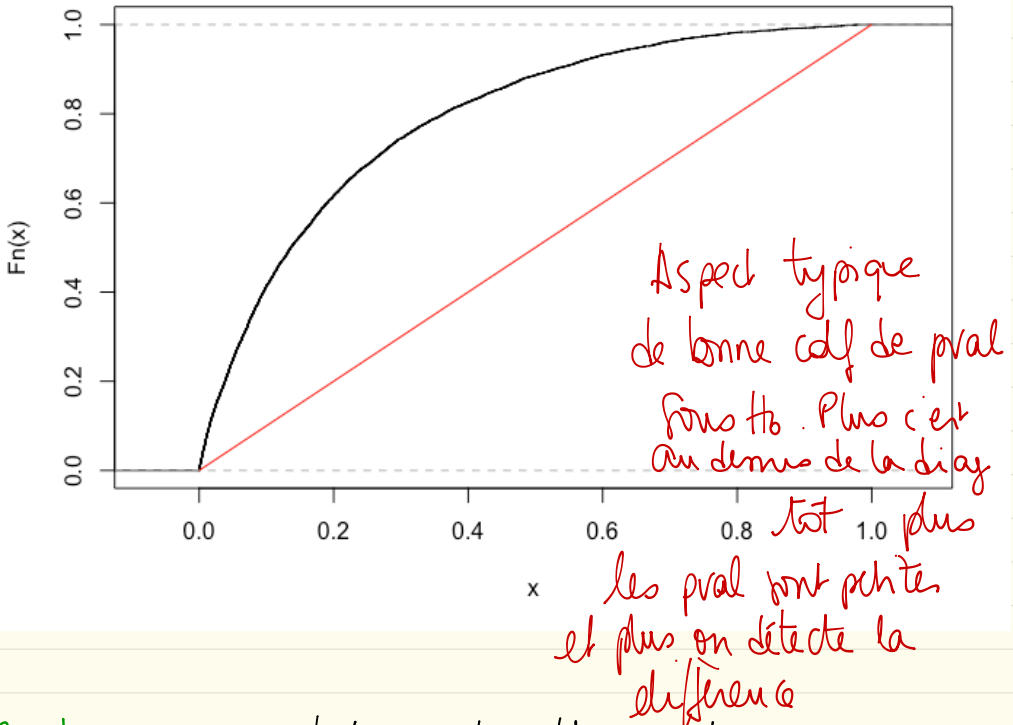
```
ks.test(X,'pnorm',sd=sqrt(2))
```

← KS voit la différence en principe.

```
for(i in 1:Nsimu)
{
  X=(2*rbinom(n,1,0.5)-1)*rexp(n,1) # l'echantillon expo et on va se demander si gaussien
  a= ks.test(X,'pnorm',sd=sqrt(2))
  pval_cont[i]=a$p.value
}
```

```
plot.ecdf(pval_cont,main='repartition des pvaleurs sous H1')
lines(c(0,1),c(0,1),col='red')
```

repartition des p valeurs sous H1



2. Les variantes (change la métrique et les quantiles mais ne change rien au final)

Cramer-von Mises utilise la distance L_2

$$n \int (\hat{F}_n(t) - F_0(t))^2 dF_0(t)$$

si besoin de la renormalisation

NB: dans le package `gofest` on les trouve
 NB2: `gof` = goodness of fit

Anderson-Darling

$$n \int \frac{(\hat{F}_n(t) - F_0(t))^2}{F_0(t)(1 - F_0(t))} dF_0(t)$$

3. Le test de KS à deux échantillons

$X_1 \dots X_n$ iid de cdf F et $Y_1 \dots Y_m$ iid de cdf G

On veut tester $H_0 : "F = G"$ contre H_1 , ce n'est pas le cas

c'est une hypothèse nulle composite

↳ il faut trouver une statistique dont la loi ne dépend pas de la valeur sous-jacente de $F = G$

$$D_{n,m} = \sup_t |\hat{F}_n(t) - \hat{G}_m(t)|$$

La loi ne dépend pas de $F = G$ si $\underline{c^0}$

↳ on rejette quand $>$ au quantile d'ordre $1 - \alpha$

$n=1000$

$X = \text{rnorm}(n)$

$Y = (2 * \text{rbinom}(n, 1, 0.5) - 1) * \text{rexp}(n, 1)$

$\text{ks.test}(X, Y)$

jouez avec n pour vous rendre compte qu'il va être beaucoup plus difficile de voir que "deux

échantillons sont \neq par rapport à

"cet échantillon ne suit pas cette loi"

NB: En pratique si les F pas c^0 , c'est que X est discret
⇒ tests du χ^2 (à revoir)

II Les tests de rang de Wilcoxon (1945)

1. Le signed-rank test

Si $X_1 \dots X_n$ iid on veut tester

H_0 : la distribution de X est symétrique autour de 0
Contre H_1 : ce n'est pas le cas.

Idee: Si la distrib^{de X} est symétrique

on peut toujours la voir comme $\text{sign}(X) |X|$

avec $\text{sign}(X)$ Rademacher (ie vaut ± 1 avec proba $1/2$)
sous H_0

Pour rendre libre de la loi de $|X|$ on transforme les $|X_i|$

en leur rang R_i . NB rang de $|X_i|$ ou rang de $F_{11}(|X_i|) \sim U_{i,i}$ (c'est pareil).
↓ c'est de $|X|$

Si $X_i \neq 0$ ps $W = \sum_{i=1}^n \text{sign}(X_i) R_i$: la loi ne dépend pas de la distrib de $|X|$

\Rightarrow accès au quantile
et rejet quand trop grand
ou trop petit

Intérêt principal = Comparaison de deux échantillons appariés

à chaque expé i , on observe X_i et Y_i , on veut savoir si $X_i > Y_i$:
 (X_i, Y_i) iid

On pose $Z = X - Y$ et on teste

H_0 : distrib de Z symétrique
 X est autant au dessus qu'en dessous de Y

H_1 : elle est décentré vers \mathbb{R}_+
 X est plus souvent au dessus de Y

revient à rejeter pour grande valeur de W unig⁺

↳ Wilcoxon-test (x, y , alternative = "greater", paired = TRUE)

2. Le rank-sum test (aussin appelé test de Mann-Whitney)

On observe X_1, \dots, X_n iid indépendants de Y_1, \dots, Y_m (et ils ne sont pas appariés!!)

$\sim X$ $\sim Y$

On suppose que $X \neq Y$ et l'on veut savoir si X est significativement plus grand que Y

Sous H_0 : $X \sim Y$, $P(X > Y) = P(Y > X)$

donc $P(X > Y) + 0.5 P(X = Y) = P(Y > X) + 0.5 P(Y = X)$

et comme la somme = 1, $P(X > Y) + 0.5 P(X = Y) = 0.5$

l'alternative est donc comprise comme H_1 : $P(X > Y) + 0.5 P(X = Y) > 0.5$
pour dire X plus "grand" que Y

Stat de test

$$U = \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}_{y_j < x_i} + 0.5 \mathbb{1}_{y_j = x_i}$$

NB si $X \sim Y$, la loi de U ne dépend pas de la loi sous-jacente

↳ on utilise son quantile pour rejeter qd U dépasse ce quantile d'ordre $1-\alpha$

↳ Wilcoxon test (X, Y , alternative = greater)

↑
par défaut pas apparié

NB: variante two sided / lower etc...
(défaut)

NB2: Ce n'est pas la même alternative que KS deux échantillons one-sided.

III Tests de normalité (cf package nortest)

On observe X_1, \dots, X_n iid de $d.f F$, on note Φ_{m, σ^2} la cdf de $N(m, \sigma^2)$

On veut tester

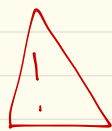
$H_0: \exists m, \sigma^2 / F = \Phi_{m, \sigma^2}$ (contre ce n'est pas le cas)

1. Lilliefors

Si on connaissait m et σ^2 on pourrait utiliser

$$D_n = \sup_t \left| F_n(t) - \Phi_{m, \sigma^2}(t) \right|$$

↳ Comme on ne les connaît pas, on pourrait faire du plug-in
= remplacer m et σ^2 par \hat{m} et $\hat{\sigma}^2$



ça change la loi de D_n (ce qui n'est pas grave
↳ il faut vérifier que ça ne dépend pas de m et σ^2 connus sous H_0 en soi)

$$T_n = \sup_t \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq t} - \Phi\left(\frac{t - \hat{m}}{\hat{\sigma}}\right) \right| \quad \text{on pose } u = \frac{t - \hat{m}}{\hat{\sigma}}$$

$$= \sup_u \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq \hat{m} + u \hat{\sigma}} - \Phi(u) \right|$$

col de $N(0,1)$

$$= \sup_u \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \leq u} - \Phi(u) \right| \quad \text{où } Y_i = \frac{X_i - \hat{m}}{\hat{\sigma}}$$

$$\text{Si } \hat{m} = \bar{X} \text{ et } \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

alors la loi de $\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ ne dépend pas de m et σ^2 si les Y_i ne sont pas \perp .

⇒ la loi de T_n sous H_0 est tabulée ⇒ le test (écrire le ...)

⇒ lillie test du package nortest en R.

2. Les variantes

On peut faire le même genre de manip sur les tests de Cameron Van Mises et Anderson (par abus de langage, on les appelle tjrs pareil)

3. Chi-deux

↳ aussi dans le package nortest

On peut tjrs faire du test du chi-deux en regroupant par classe. (m variable pour KS à 1 éch)

Si m et σ^2 pas connus → Test avec est⁰ de paramètres 1 du chi-deux (+/- vrai con distri.) pas discrète

(Vous avez dû le voir l'année précédente)

↳ aussi dans le package nortest

4. Shapiro et Wilk

↳ shapiro.test connu pour être le + puissant

aller voir la formule si vous êtes curieux.

