

### Bernard Billoud Morphogenesis of Macro-Algae - UMR8227 Bernard.Billoud@sb-roscoff.fr



June 1, 2021

# Long is the day

### Introduction: what it is all about

- Setting up a test
- 2 The normal law
  - From discrete to continuous
  - Properties
  - Cumulative Density Function
- 3 Sampling
  - Normal
  - not Normal
  - Sample size and convergence

### 4 Testing

- one sample against a normal population
- two samples
- paired samples
- variance
- more than two samples
- normality
- non-normal variables
- small samples of unknown distribution
- counts in classes
- linear dependance
- combined effects
- non-linear dependance
- 5 Conclusion
  - Take-home graph

#### Setting up a test

# The father of all tests: the Gauss test



What proportion of all "usually-washed" T-shirts are whiter than the one washed with the new OMO?



# The T-Shirt test (enhanced version)

Threshold Choose a confidence level, *e.g.*  $\alpha$  = 0.05;

Null hypothesis  $H_0$  Let us suppose that the new OMO is *not* better than my usual powder;

Variable distribution Knowing the distribution of W, the whiteness of my  $N_{tot} = 216$  regular T-Shirts;

Test value Knowing *W*<sub>O</sub> the whiteness of the T-Shirt I washed using the new OMO;

p-value The probability to find a regular T-Shirt as white as (or whiter than) the OMO T-Shirt is:

$$P(W \ge W_O) = \frac{N(W \ge W_O)}{N_{tot}}$$



As  $P = 0.06 \ge \alpha$ , it can be expected that by chance, a T-Shirt picked at random among the regular ones would be at least as white as the one I washed with the new OMO.

So:

I cannot assert that the new OMO washes whiter than my usual powder.

B.B. (MMA / SBR)

## Code for the T-Shirt test

```
# T-Shirt whiteness distribution -- version 1
tsw <- read.table("RegularWashing.csv",header=T)</pre>
B = (min(tsw$whiteness)-0.5):(max(tsw$whiteness)+0.5)
h <- hist(tsw$whiteness.breaks=B)</pre>
W = hsmids
N = hscounts
# T-Shirt whiteness distribution -- version 2
\mathbb{N} = \mathbf{c}(1, 1, 2, 2, 3, 4, 5, 7, 8, 9, 10, 11, 12, 13, 13, 14, 13, 13, 12, 11, 10, 9, 8, 7, 5, 4, 3, 2, 2, 1, 1)
W = seq(0.length(N)-1)
# Whiteness of the OMO-washed T-Shirt and number of whiter T-Shirts
W0 = 25
listOK = N[W > = WO]
nOK = sum(listOK)
# Test
Ntot = sum(N)
cat("\n",Ntot," T-Shirts, in which ",nOK," have whiteness >= ",WO,": ",sep="")
print(listOK)
pv = n0K / Ntot
cat("p-value =",pv,"\n")
```

# Counting and computing

A simple model of allocation in classes: binomial (discrete)



- The value X (position of one bead) depends on multiple (*n* rows) "choices" (L/R) with probability *p* (= 1/2)
- Distribution among n + 1 values (bucket  $0 \le k \le n$ ):

$$P(X = k) = B_{n,p}(k) = \frac{n!}{k!(n-k)!}p^k(1-p)^{n-k}$$

When *n* increases:



• value range extends:  $X \rightarrow$  continuous

• 
$$B_{n,p} \to \mathcal{N}_{np,np(1-p)}$$

• Density of probability:

$$P(x) = N_{\mu,\sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

# Interlude: installing a package

R Studio allows any user to – locally – install a package

#### install.packages("Rlab")

File Edit Code View Plots Session Build Debug Profile	Tools Help	
🔍 🔍 🗸 😪 🚭 🗸 📑 📑 📄 🍙 Go to file/function 👘 🔡 🗸 Addin	Install Packages	
Console Terminal × Jobs ×	Check for Package Updates     Connections	
/home/umr7139/mma/billoud/ 🗇	Version Control Install R packages	
	Shell	Install Packages
the set of a set of the set of the	Terminal	Install from: (2) Configuring Repositories
	A store	Repository (CRAN) V
A REAL PROPERTY AND AN AN AND AND AN	Addins •	Packages (separate multiple with space or comma):
	Keyboard Shortcuts Help Shift+Alt+K Modify Keyboard Shortcuts	Rlab r. Rlabkey Rlabkey Rlabkey 39/mma/billoud/R/x86_64-pc-linux-gnu-librar∽
Contract The second second second second	Project Options	rlang rlas idencies
	Global Options	
the second s		Install Cancel

## Code for the Galton board

```
library(Rlab) # Rlab is used for Bernouilli assay (rbern)
```

```
# Draw 0 or 1 with proba 0.5 (= 1 pod of the Galton board)
proba = 0.5
d = rbern(1.proba)
cat("\nResult of one Bernouilli assay with p = ",proba,": ",d,"\n",sep="")
# Assuming 0 = L; 1 = R, draw 14 times 0 or 1 with proba 0.5 -> bucket number for 1 bead
nlin = 14
x = sum(rbern(nlin, proba))
cat("\nSum of n = ".nlin." Bernouilli assays with p = ".proba.": ".x."\n".sep="")
# Looping
nb = 10000 # Launch 10000 beads
X < - c()
for (i in 1:nb) {
  X \le c(X.sum(rbern(nlin.proba)))
}
m = mean(X) \# expected: 14 \times 0.5 = 7
s = sd(X) # expected: \sqrt{(14 \times 0.5 \times 0.5)} = 1.87
cat(nb," sums of ",nlin," assays: mean = ",m,"; s.d. = ",s,"\n",sep="")
cat("Expected: mean = ",nlin*proba,"; s.d. = ",sqrt(nlin*proba*(1-proba)),"\n",sep="")
# Show
hist(X,breaks=-0.5:(nlin+.5),xlim=c(0,nlin),col="orange",
   main=paste("Binomial law: n = ".nlin."; p = ".proba." (pop = ".nb.")".sep=""))
```

## Improve efficiency and converge to the Normal law

```
cat("\nUsing rbinom:\n")
```

```
proba = 0.5 ; nlin = 14 ; nb = 10000
# Use the appropriate R function
X = rbinom(nb, nlin, proba)
m = mean(X) \# expected: 14 \times 0.5 = 7
s = sd(X) # expected: \sqrt{(14 \times 0.5 \times 0.5)} = 1.87
cat(nb," sums of ",nlin," assays: mean = ",m,"; s.d. = ",s,"\n",sep="")
cat("Expected: mean = ",nlin*proba,"; s.d. = ",sqrt(nlin*proba*(1-proba)),"\n",sep="")
cat("\nIncreasing n:\n")
# Increase nlin (try also with nb = 100000)
nlin = 436
X = rbinom(nb. nlin. proba)
m = mean(X) \# expected: 436 \times 0.5 = 218
s = sd(X)
          # expected: \sqrt{(436 \times 0.5 \times 0.5)} = 10.44
hist(X.breaks=-0.5:(nlin+.5).col="vellow".freg=F.
   xlim=c(nlin*proba-5*sqrt(nlin*proba*(1-proba)).nlin*proba+5*sqrt(nlin*proba*(1-proba))).
   main=paste("Binomial law: n = ",nlin,"; p = ",proba," (pop = ",nb,")",sep=""))
# Compare with Normal
valist = 0·nlin
normd = dnorm(valist,m,s)
points(valist.normd.type="1".lwd=3.col="blue")
cat(nb." sums of ".nlin." assays: mean = ".m.": s.d. = ".s."\n".sep="")
cat("Expected: mean = ",nlin*proba,"; s.d. = ",sqrt(nlin*proba*(1-proba)),"\n",sep="")
```

## The *real* Gauss test = Z-test



What proportion of all "usually-washed" T-shirts are whiter than the one washed with the new OMO?





## Code for the Z-test

# Load the library for Gauss.test
library(compositions)

```
# T-Shirt whiteness distribution
W = 0:30
N = c(1,1,2,2,3,4,5,7,8,9,10,11,12,13,13,14,13,13,12,11,10,9,8,7,5,4,3,2,2,1,1)
Ntot = sum(N)
```

```
# Parameters for the norml law
meanW = sum(W*N) / Ntot
sdW = sart(sum(W*W*N) / Ntot - meanW**2)
cat("\nWhiteness: mean =".meanW."s.d. =".sdW."\n")
# Discrete: WO = 25, test for W \ge WO; Continuous: test for W \ge WO, so:
WO = 24.5
# Test
zt = Gauss.test(W0.mean=meanW.sd=sdW.alternative="greater")
# print(zt)
pv = zt p. value
cat("p-value for the z-test (W>",WO,"): ",pv,"\n",sep="")
# Graphic representation
valist <- seg(meanW-5*sdW.meanW+5*sdW.length.out=100)</pre>
dens <- dnorm(valist,m=meanW,s=sdW)</pre>
plot(x=valist,y=dens,type="1",col="darkmagenta",lwd=3,xlab="x",ylab="Probability density at x",
     main=paste("Distribution N(",meanW,",",round(sdW,3),")",sep=""))
abline(v=W0,lwd=3,col="red")
```

# Representative curve

### With different values for $\mu$ and $\sigma$



## Representative curve

### With different values for $\mu$ and $\sigma$



### All curves have the same shape

#### Properties

## Parameters $\mu$ and $\sigma$

Normal law: changing  $\mu$ 





A somewhat useful formula:  $P_2(x_2) = P_1(x_1) = P_1(x_2 - (\mu_2 - \mu_1))$  A completely useless formula:  $P_2(x_2) = \frac{\sigma_1}{\sigma_2} P_1(x_1) = \frac{\sigma_1}{\sigma_2} P_1\left(\frac{\sigma_1}{\sigma_2} x_2 + \mu\left(1 - \frac{\sigma_1}{\sigma_2}\right)\right)$ 

# Mapping to the standard normal distribution



A really useful formula:  $P_{\mu,\sigma}(x) = \frac{1}{\sigma} P_{0,1}\left(\frac{x-\mu}{\sigma}\right)$ 

$$Or: X \mapsto N_{\mu,\sigma} \Leftrightarrow X \mapsto \mu + \sigma \times N_{0,1}$$

# Density of probability and area

Actually, we are interested by the area under the curve



### F is versatile



### Code for various cases

```
# Suppose we have a population P. for which we know \mu and \sigma
mu = 218
sigma = 10
cat("H0: x drawn from population: N(\mu,\sigma) with \mu = ".mu," and \sigma = ".sigma."\n".sep="")
valist = c(198, 203, 208, 218, 228, 233, 238)
# pnorm directly gives the probability that x < u
cat("\nLeft test:\n")
for ( u in valist ) {
  cat("Under HO, P(x<",u,") = F(",u,") = ",pnorm(u,m=mu,s=sigma),"\n",sep="")
# pnorm indirectly gives the probability that u < x
cat("\nRight test:\n")
for ( u in valist ) {
  cat("Under H0, P(",u,"<x) = 1 - F(",u,") = ",1-pnorm(u,m=mu,s=sigma),"\n",sep="")
}
# pnorm indirectly gives the probability that u1 < x < u2
cat("\nInside interval:\n")
for ( i in 1:(length(valist)-1) ) {
  u1 = valist[i]
  for ( j in (i+1):length(valist) ) {
    u2 = valist[j]
    cat("Under H0, P(",u1,"<x<",u2,") = F(",u2,") - F(",u1,") = ",</pre>
                          pnorm(u2.m=mu.s=sigma)-pnorm(u1.m=mu.s=sigma)."\n".sep="")
```

## Some rules of thumb



#### Cumulative Density Function

# Example of application

The size of one year old *Cerastoderma edule* follows a normal law with  $\mu = 11.45$ mm;  $\sigma = 2.86$ mm<sup>1</sup>. We find a shell which seems to be a young *C. edule*, and we measure its size: S = 17mm. Is our specimen too large to be a young *C. edule*?







B.B. (MMA / SBR)

# Sampling in the normal law



240

260

#### not Normal

# Let us define a new law: the RC law



## Small samples

### Same $\mu$ = 50.05 and $\sigma$ = 22.31



B.B. (MMA / SBR)

# Convergence of $\bar{X}$



Increasing sample size : Central Limit Theorem

$$\lim_{n \to +\infty} \bar{X} = N_{\mu,\sigma/\sqrt{n}}$$
$$\hat{\mu} = \bar{x}$$

# Convergence of $\bar{S}$



Whatever the sample size :

$$\hat{\sigma} = \sqrt{\frac{n}{n-1}s}$$

#### Sample size and convergence

# Confidence Interval

Question<sup>2</sup>: How credible is a parameter estimate?

Example: dry weight of *Ophiothrix fragilis* females



# Measures nOf = 161 mOf = 696.00 sOf = 134.23 # Estimates for the population muOf = mOf siOf = sOf \* sqrt(nOf/(nOf-1)) n = 161, m = 696.00 mg, s = 134.23 mg.

Estimates for the population:

$$\hat{\mu} = m$$
  $\hat{\sigma} = s\sqrt{\frac{n}{n-1}} = 134.65 \text{ mg}$ 

Confidence interval at 95%, *i.e.*  $\alpha = 0.05$ 

$$\hat{\mu} \pm u_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

# Confidence interval
alpha = 0.05
rnglo = qnorm(alpha/2,mu0f,si0f / sqrt(n0f))
rnghi = qnorm(1-alpha/2,mu0f,si0f / sqrt(n0f))
cat("Confidence interval at 95%: [",rnglo,":",rnghi,"]\n")

<sup>2</sup> Inspired by A. Lefebvre, D. Davoult, F. Gentil, M. A. Janquin (1999) *Hydrobiologia* 414:25–34.

B.B. (MMA / SBR)

Stats with F

Sampling

Sample size and convergence

# Laws and sample size: the sample, that's simple



B.B. (MMA / SBR)

### Student's Law

valist <- seq(-5,5,length.out=200)
d = 2 # degree of freedom
plot(valist,dt(valist,df=d),type="1")</pre>



### Student's law is like a Normal law with heavy tails (for small d.f.)

3e-04

# Example of application

The biomass of green algae in an estuary is known to follow a normal law with  $\mu = 3906 \text{ g/m}^2$ ;  $\sigma = 599 \text{ g/m}^2$ .

We suspect a pollution to increase this mass.

We perform 10 measures and obtain the data in file MassAlg01.tab Should we confirm that the mass has been increased ?

```
# Under HO, the sample follows the normal law of the population:
                                                                       Density
cat("HO : mass does not exceed the usual norm.\n")
                                                                         te-04
muA = 3906; sigmaA = 599
```

```
# Observed data:
massalg <- read.table("MassAlg01.tab",header=F)</pre>
nbmes = nrow(massalg)
                                                                         00+90
measM = mean(massalg$V1)
                                                                              3000
cat("Mean mass on", nbmes, "measures: ", round(measM,2), "\n")
                                                                             N = 10
                                                                                     Bandwidth = 220.7
plot(density(massalg$V1).main="Algal mass")  # Rather nice, but not so rigorous
```

```
# Under H0, (\bar{X}-\mu)/(\sigma/_{1}h) \mapsto t(n-1 df)
alpha = 0.05
pv = 1-pt((measM-muA)/(sigmaA/sqrt(nbmes)),df=nbmes-1)
cat("p-value =",pv,"; ")
if(pv<alpha) {</pre>
  cat("we can reject HO")
} else {
  cat("we cannot reject HO")
}
cat(" with a type I error risk \alpha =".alpha."\n\n")
```

5000

6000

Algal mass

observed  $N(\bar{x}.s)$ N(μ,σ)

4000

## Two changes

### What if $\sigma$ is not known?

If we do not know  $\sigma$ , we replace it by our estimate:  $\hat{\sigma} = s\sqrt{\frac{n}{n-1}}$ Warning ! R computes sd with (n-1) so, it is in fact  $\hat{\sigma}$ !

```
cat("The standard deviation for the sample is : s =",
    sd(massalg$V1)*sqrt((nbmes-1)/nbmes),"\n")
estimsig = sd(massalg$V1)
cat("The estimated population s.d. is", estimsig,"\n")
```

```
# recompute p-value using this estimate
pv = 1-pt((measM-muA)/(estimsig/sqrt(nbmes)),df=nbmes-1)
```

```
cat("Using this estimate, we find p-value =",pv,"\n")
```

# What if we are lazy?

Use The R implementation of the t test!

```
cat("\n[R] : function t.test\n")
tt <- t.test(massalg$V1,mu=muA,alternative="greater")
print(tt)</pre>
```

# The real experimental situation

In real experiments, we usually have two (or more) samples to compare



# Distribution and field of application

- $\overline{x_A x_B}$  and  $\overline{x_A} \overline{x_B}$  have the same mean and s.d., so how to choose?
  - $\overline{x_A} \overline{x_B}$  can be computed for unpaired samples, even if  $n_A \neq n_B$

$$\overline{x_A} - \overline{x_B} \mapsto (\mu_A - \mu_B) + \sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}} \times t_{n_A - 1 + n_B - 1}$$

Usually we do not know  $\sigma_A$  and  $\sigma_B$ :

$$\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B} \approx (n_A s_A^2 + n_B s_B^2) \times \frac{1/n_A + 1/n_B}{n_A + n_B - 2}$$

•  $\overline{x_A - x_B}$  is good for pairwise differences (thus  $n_A = n_B = n$ )

$$\overline{x_A - x_B} \mapsto (\mu_A - \mu_B) + \sqrt{\frac{\sigma_A^2 + \sigma_B^2}{n}} \times t_{n-1}$$

Usually we do not know  $\sigma_A$  and  $\sigma_B$ :

$$\frac{\sigma_A^2 + \sigma_B^2}{n} \approx \frac{s_A^2 + s_B^2}{n - 1}$$



<sup>3</sup> A. Fort, C. Mannion, J.M. Fariñas-Franco, R. Sulpice (2020) Sci. Tot. Envir. 698 134337

B.B. (MMA / SBR)

Stats with R

### Why experimentalists hate statisticians



We cannot change anything to  $\mu_A$ ,  $\mu_B$ ,  $\sigma_{x_A}$ ,  $\sigma_{x_B}$ 

The only parameter we can tune:

$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$$

Increasing *n* decreases dispersion of means

 $\Leftrightarrow$ The larger *n*, the best we see (small) effects

More bench required!

# Improve your power by using the pairing information

Example: Genotypic diversity in Agarophyton chilense<sup>4</sup>

Question: Effect on growth rate (supposed normally-distributed)?

```
sgrdata <- read.table("SGR_Chaica.tab",header=T)</pre>
attach(sqrdata)
                                                                   Genetic diversity and growth rate
summary(sgrdata[,2:3])
                                                           Specific Growth Rate (%/day
# default: unpaired
cat("\nNot taking pairs into account:\n")
cat(" Difference of means =",mean(T1G)-mean(T4G),"\n")
t.test(T1G,T4G,var.equal=T)
# actually, measures on the same genotype => paired
cat("\nTaking pairs (Genotypes) into account:\n")
cat(" Mean difference =".mean(T1G-T4G)."\n")
t.test(T1G,T4G,var.equal=T,paired=T)
# draw
boxplot(T1G,T4G,col=c("purple","orange"),
  main="Genetic diversity and growth rate",
                                                                       1 aenome
  xlab="Diversity",ylab="Specific Growth Rate (%/day)",
                                                                               Diversity
  names=c("1 genome","4 genomes") )
```

### See the difference? (Hint: pay attention to the df)

<sup>4</sup> Data dishonestly extracted from S. Usandizaga, A.H. Buschmann, C. Camus, J.L. Kappes, S. Arnaud-Haond, S. Mauger, M. Valero and M.L. Guillemin (2019) Evol. Appl. 00:1-13.

B.B. (MMA / SBR)

4 aenomes

# How do we know if variances are equal?

Fisher:

- One population with variance  $\sigma^2$
- Two samples with variance  $s_A^2$  and  $s_B^2$
- Compute the law for  $\frac{s_A^2}{s_B^2}$


### Example

```
muA = 60 ; muB = 55
sigmA = 3; sigmB = 6
nA = 12 : nB = 15
listA <- rnorm(nA, m=muA, s=sigmA)</pre>
listB <- rnorm(nB, m=muB, s=sigmB)
mA = mean(listA)
mB = mean(listB)
sA = sd(listA)
sB = sd(listB)
cat("sA =".sA.": sB =".sB.": F =".(sA/sB)**2)
# Bilateral test
pv = 2 * pf((sA/sB)**2, df1=nA-1, df2=nB-1)
cat(" p-value =",pv,"\n")
# Too easy
cat("\n[R] : function var.test\n")
var.test(listA, listB) # test F, same p-value as var.test(listB, listA)
```

# So, if variances are not equal?

• 
$$\hat{\sigma}_A$$
 and  $\hat{\sigma}_B \rightarrow \sigma^2 = \frac{\hat{\sigma}_A^2}{n_A} + \frac{\hat{\sigma}_B^2}{n_B}$ 

• change the degree of freedom to something complicated  $\nu < n_A + n_B - 2$  and  $\nu \rightarrow n_A + n_B - 2$  when  $s_A \rightarrow s_B$ 



### Comparing more than two samples

Example: effect of (S)-roscovitine on metabolic parameters in rats<sup>5</sup>. Body temperature for 4 conditions,



Untreated	Vehicle	(S)-Roscovitin	(S)Roscovitin 10 mg/kg/h
	(Control)	Pre-occlusion	Post-occlusion
38.9	38.8	37.9	38.7
38.7	38.8	38.5	38.6
38.6	38.3	38.2	38.5
38.8	38.5	38.3	39.0
38.9	38.1	38.3	38.9
38.8	38.8	38.7	38.3
38.9	38.5	38.3	38.6
38.8	38.4	38.3	38.4
38.9	38.6	38.2	38.6
38.7	38.4	38.1	38.7
38.8	38.6	38.3	38.5
38.9	38.6		38.4
	38.4		

with  $11 \le n \le 13$ 

### First idea: One test for each pair of conditions

```
# Read data
neurokin <- read.table("Neurokin_01.tab",header=T)
attach(neurokin)</pre>
```

# Plot
boxplot(neurokin)

#perform all 6 tests

- t.test(Untreated,Vehicle)
- t.test(Untreated,SRoscoPre)
- t.test(Untreated,SRoscoPost10)
- t.test(Vehicle,SRoscoPre)
- t.test(Vehicle,SRoscoPost10)
- t.test(SRoscoPre,SRoscoPost10)

<sup>5</sup> Data reconstructed from

B. Menn, S. Bach, T. Blevins, M. Campbell, L. Meijer, S. Timsit (2010) PLoS ONE 5(8) e12117.

B.B. (MMA / SBR)

### Interlude: reshaping data

### Some functions require an appropriate data format.

"wid	le"	format

Untreated	Vehicle	SRoscoPre	SRoscoPost10
38.9	38.8	37.9	38.7
38.7	38.8	38.5	38.6
38.6	38.3	38.2	38.5
38.8	38.5	38.3	39.0
38.9	38.1	38.3	38.9
38.8	38.8	38.7	38.3
38.9	38.5	38.3	38.6
38.8	38.4	38.3	38.4
38.9	38.6	38.2	38.6
38.7	38.4	38.1	38.7
38.8	38.6	38.3	38.5
38.9	38.6		38.4
	38.4		

### From wide to long:

```
neuroCT <- stack(neurokin)
neuroCT <- neuroCT[!is.na(neuroCT$values),c(2,1)]
colnames(neuroCT) <- c("Condition", "Temperature")</pre>
```

#### or:

"long" format

	C 111	
	Condition	Iemperature
1	Untreated	38.9
2	Untreated	38.7
3	Untreated	38.6
10	Untreated	38.7
11	Untreated	38.8
12	Untreated	38.9
14	Vehicle	38.8
15	Vehicle	38.8
16	Vehicle	38.3
24	Vehicle	38.6
25	Vehicle	38.6
26	Vehicle	38.4
27	SRoscoPre	37.9
28	SRoscoPre	38.5
29	SRoscoPre	38.2
35	SRoscoPre	38.2
36	SRoscoPre	38.1
37	SRoscoPre	38.3
40	SRoscoPost10	38.7
41	SRoscoPost10	38.6
42	SRoscoPost10	38.5
49	SRoscoPost10	38.7
50	SRoscoPost10	38.5
51	SRoscoPost10	38.4

### Performing many t-test is pleasurable, but...

*H*<sub>0</sub> should be: "All samples are drawn out of the same population" ⇒ same σ for all conditions but : each t-test computes its own ô

*p* value over- or under-estimated, depending on  $\hat{\sigma}_{ij}$  vs  $\sigma$ 

*p* value says how probable are the data under H<sub>0</sub>
 but : each t-test computes how probable is its own subset of two samples

- Risk  $\alpha_p = 0.05$  : error on 1 pairwise test
- ►  $\Rightarrow$   $P_p = 1 \alpha_p = 0.95$  for each pairwise test to be correct
- ► 6 simultaneously correct tests:  $P_M = P_p^6 = (1 \alpha_p)^6 \approx 0.735$
- ► ⇒ we accept a risk  $\alpha_M \approx 0.265$  of error for the whole analysis... do we?

Rule: pairwise *p* values must be adjusted for multiple testing

```
#perform all 6 tests
# Read data
                                                         pairwise.t.test(stneu$Temperature,
neurokin <- read.table("Neurokin_01.tab",</pre>
                                                                         q=stneu$Condition)
                                     header=T)
                                                            Pairwise comparisons using t tests with pooled SD
attach(neurokin)
                                                       data: stneu$Temperature and stneu$Condition
# Re-format data
                                                                    Untreated Vehicle SRoscoPre
stneu <- stack(neurokin)</pre>
                                                       Vehicle
                                                                    0.00152
stneu <- stneu[!is.na(stneu$values),c(2,1)]</pre>
                                                       SRoscoPre
                                                                    1.3e-07
                                                                              0.00808 -
                                                       SRoscoPost10 0.01697
colnames(stneu) <- c("Condition"."Temperature")</pre>
                                                                              0.30506 0.00083
                                                       P value adjustment method: holm
```

```
B.B. (MMA / SBR)
```

## Analysis Of Variance

Reminder: under  $H_0$ ,

The variation of mean between classes is not higher than the mean of variation within each class

$$N = \sum_{j=1}^{k} n_j$$
  $M = \frac{1}{N} \sum_{j=1}^{k} \sum_{i=1}^{n_j} x_{i,j}$ 

Mean for each class:

 $m_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{i,j}$ 

Variance within each class:

$$s_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (x_{i,j} - m_j)^2$$

Variance due to Factor:

Intrinsic variance:

$$S_F^2 = \frac{1}{k-1} \sum_{j=1}^{\kappa} n_j (m_j - M)^2$$

$$S_l^2 = \frac{1}{N-k} \sum_{j=1}^k (n_j - 1) s_j^2$$

Hence the question: is  $S_F^2$  significantly larger than  $S_I^2$ ? Hence the question: how do we compare variances? Hence the answer: with a Fisher test.

B.B. (MMA / SBR)

### Let's write some code!

```
neurokin <- read.table("Neurokin_01.tab",header=T)
# convert to long format
stneu <- stack(neurokin)
stneu <- stneu[!is.na(stneu$values),c(2,1)]
colnames(stneu) <- c("Condition", "Temperature")
attach(stneu)</pre>
```

```
# basic numbers
N = nrow(stneu)
co = levels(Condition)
k = length(co)
M = mean(Temperature)
# loop to sum into S2F and S2I
S2F = 0
S2I = 0
for (j in 1:k) {
  xj <- stneu[stneu[,1]==co[j],2]</pre>
  n_i = length(x_i)
  S2F = S2F + nj / (k-1) * (mean(xj) - M) * 2
  S2I = S2I + (nj-1) / (N-k) * var(xj)
# perform a Fisher test
Ffi = S2F/S2I
cat("k =",k," N =",N," k-1 =",k-1," N-k =",N-k,"\n")
cat("S2b =",S2F," S2w =",S2I," Ffi =",Ffi,"\n")
pv <- 1 - pf(Ffi, df1=k-1, df2=N-k)
cat("p-value for AnOVa test: p = ", pv, " \setminus n \setminus n")
```



# Let's avoid writing so much code!

# use the R function aov
summary(aov(Temperature ~ Condition))

### What we learn from a one-way AnOVa

A typical one-way AnOVa test:

- Choose α
- *H*<sub>0</sub>: All samples come from the same population distibuted according to the same Normal law:

 $\mu_1=\mu_2=\ldots=\mu_k,\,\sigma_1=\sigma_2=\ldots=\sigma_k$ 

- Compute  $F = S_F^2/S_I^2$ , then p
  - if *p* > *α*, then *H*<sub>0</sub> cannot be rejected
     ⇒ No effect can be demonstrated for the factor
  - if p < α, then H<sub>0</sub> must be rejected
     ⇒ (provided normality and homeosedasticity hold)
     The factor affects μ, at least between two classes

One-way AnOVa test does not tell us which classes differ!

To discover where is/are the difference(s):

```
neurAOV <- aov(Temperature~Condition)
TukeyHSD(neurAOV)</pre>
```

The Tukey's "*post-hoc*" test is similar to the pairwise *t*-test with *p*-value correction = "Honnest Significant Difference" (  $\approx$  Bonferroni).

#### Testing normality

### Testing normality

Is a given sample likely to have been drawn out of a Normal law ?



B.B. (MMA / SBR)

Stats with R

### Testing normality

Is a given sample likely to have been drawn out of a Normal law ?

Indices: for Normal law, skewness = 0; kurtosis = 3 Example: Gene expression level







### Testing normality

Is a given sample likely to have been drawn out of a Normal law ? Example: convergence of Student's law to the Normal when  $\nu$  increases

```
listn <- rnorm(n=2000)
for ( d in c(1,5,30,100) ) {
    listt <- rt(n=2000,df=d)
    qqplot(listn,listt,xlim=c(-5,5),ylim=c(-5,5),pch=16,col=rgb(0.8,0.1,1,0.2),
        main="Convergence of t to N",
        ylab=substitute(paste("Sample of t,",nu,"=",n,"; n = 2000"),list(n=d)),
        xlab="Sample of Normal law, n=2000",
        }
    abline(0,1,col="blue")
    st = shapiro.test(listt)
    cat("v = ",d,"p-value for Shapiro-Wilk test =",st$p.value,"\n")
}</pre>
```



### Warning

### For nearly normal laws, Shapiro-Wilk test result depends on sample size

```
# create functions for Student law
# ... with \nu = 20
t20 <- function(n) {rt(n,df=20)}
# ... with \nu = 50
t50 <- function(n) {rt(n,df=50)}
```

```
# Change sample size
for(ss in c(10,30,100,1000,5000)){
```

```
cat(sprintf(" %4d ",ss))
```

```
for(law in c(rlnorm,t20,t50,rnorm)) {
    pv<-c()
    for(i in seq(1,2000)){
        x<-law(n=ss)
        pv<-c(pv,shapiro.test(x)$p.value)
    }
    # Proportion of normality rejection</pre>
```

```
cat(sprintf(" %.3f ",
```

```
length(pv[pv<0.05])/length(pv)))</pre>
```

**cat**("\n")

Proportion of normality rejection

sample	log N	Student		Normal
size n	iog-in	$\nu = 20$	$\nu = 50$	Norman
10	0.594	0.058	0.053	0.043
30	0.995	0.081	0.064	0.043
100	1.000	0.111	0.077	0.043
1000	1.000	0.400	0.117	0.046
5000	1.000	0.942	0.260	0.036
Expected	$\leftarrow$ Reject $\rightarrow$		← Ac	ccept →
Actual	≈OK	depend	ds on <i>n</i>	OK

### Poisson law for fishermen

- A fisherman can catch 0 or 1 fish every minute.
- The probability to catch 1 fish during 1 minute is  $\beta = 1/12$  (he catches on average 1 fish every 12 minutes).
- We observe him for  $\Delta t = 30$  min.

What is the probability he catches k = 3 fishes?

Poisson law of parameter  $\lambda = \beta \Delta t$ :  $P_{\lambda}(n = k) = \frac{\lambda^k}{k!} e^{-\lambda}$  $P_{\lambda}: \mu = \lambda; \sigma^2 = \lambda$ 

λ = βΔt = 30/12 = 2.5 fishes during 30 min;  $P_{2.5}(n = 3) = 2.5^3/3! \times e^{-2.5} \approx 0.214$ 

# Poisson probability of 3 with \u03c0 = 30/12
dpois(3, lambda = 30/12)
# Plot
k = 0:9
barplot(dpois(k, lambda = 30/12), col = "salmon",
 main = "Poisson Law for fisherman",
 xlab = "k", ylab = "P(n=k)"
)



### Poisson law for molecular biologists

- GAATTC can occur 0 or 1 time at any position in DNA.
- The probability that a 6-mer is GAATTC is  $\beta \approx 1/4^6 = 1/4096$  (GAATTC occurs on average 1 time every 4096 bases).
- We observe a sequence of length l = 10 kb.

What is the probability we find *k* = 3 occurrences of GAATTC in *l*?

Poisson law of parameter  $\lambda = \beta l$ :  $P_{\lambda}(n = k) = \frac{\lambda^k}{k!}e^{-\lambda}$  $P_{\lambda}: \mu = \lambda; \sigma^2 = \lambda$ 

$$\lambda = \beta \ l = 10000/4096 \approx 2.4414 \text{ cuts in } 10 \text{ kb};$$
  
 $P_{2.4414}(n = 3) = 2.4414^3/3! \times e^{-2.4414} \approx 0.211$ 





### Poisson law for NGS

- 1 read can match on 0 or 1 expressed mRNA (kind of wishful thinking).
- The probability that this read matches mRNA *m* is  $\beta = l_m n_m / \Sigma l_i n_i$  (Proportion of nucleotides in *m* among nucleotides in transcriptome).
- We observe a pool of  $N_r = 10$  millions reads.

Assume

- 5000 expressed genes with average: length = 2kb, number = 800  $\Rightarrow \Sigma l_i n_i = 5000 \times 2000 \times 800 = 8 \times 10^9$  nucleotides.
- messenger *m* is  $l_m = 2.5$ kb long, and is expressed  $n_m = 600$ times No expression level variability  $\Rightarrow$  technical replicates  $\Rightarrow \beta = 1.875 \times 10^{-4} \Rightarrow \lambda = \beta N_r = 1875$  read counts.

```
# Plot
k = 1700:2050
barplot(dpois(k, lambda = 1875), col = "forestgreen",
    main = "Poisson Law for NGS",
    xlab = "k", ylab = "P(n=k)"
) # ... hmm, this looks like a normal law!
```

```
# Compare two technical replicates: same \lambda?
poisson.test(c(1948,1843)) # does not know \lambda=1875
```



### Wilcoxon W test: principle

Consider two independent samples selected from populations...

### with the same distribution

 $\begin{array}{rll} \text{sample A} & \text{sample B} \\ & 54 & 90 & 56 & 22 & 73 & 49 & 79 & 85 & 40 & 20 \\ & 26 & 17 & 25 & 50 & 67 & 12 & 72 & 10 & 28 & 53 \\ \hline & \text{Merge and order} \\ & 10 & 12 & 17 & 20 & 22 & 25 & 26 & 28 & 40 & 49 \\ & 50 & 53 & 54 & 56 & 67 & 72 & 73 & 79 & 85 & 90 \\ \hline & \text{Forget values} \\ & 12 & 3 & 45 & 67 & 89 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 \\ \hline & \text{Sum ranks} \\ & W_A &= \sum r_A - \sum_{i=1}^{n_A} i = 96 - 45 = 51 \\ & W_B &= \sum r_B - \sum_{i=1}^{n_B} i = 114 - 66 = 48 \end{array}$ 

with different distributions

Under  $H_0$  = samples A and B are issued from the same law:  $W \mapsto \text{Wilcoxon} - \text{Mann} - \text{Whitney law} \left( E = \frac{n_A n_B}{2}; V = \frac{n_A n_B (n_A + n_B + 1)}{12} \right)$  $E = 49.5; V = 173.25; \sqrt{V} = 13.16$ 

$$p = 0.94$$
 |  $p = 4.24 \times 10^{-3}$ 

### Example: shaking oysters

### Question:

Does a mechanical stress change Oyster's survival to microbial infection?

Method<sup>6</sup>:

Infect, then leave alone or shake 5 min on Day 3; count mortality on Day 9.

But: Only 4 replicates, no clue about how survival distibutes.

```
# Read data
od <- read.table("OysterDeath.tab",header=T)
attach(od)</pre>
```

```
#perform the test
wilcox.test(Mortality~Cond)
```

```
nA = length(Mortality[Cond=="ChalNoStress"])
nB = length(Mortality[Cond=="Chal5mStress"])
cat("E(W) =",nA*nB/2,"\n")
```



<sup>6</sup> Inspired by A. Lacoste, F. Jalabert, S.K. Malham, A. Cueff, S.A. Poulet (2001) Appl. Envir. Microbiol. 67:5 2304–2309.

### Testing numbers of individuals - Scheme

### Population ( $\infty$ ) Sample N=300

		Class	Droha	Exported	Sampla	· <sup>2</sup>
		Class	Proba	Expected	Sample	Xi
		i	Pi	$e_i = Np_i$	si	$\frac{(s_i - e_i)^2}{e_i}$
			0.052	15.6	18	0.3692
200		$\bigcirc$	0.097	29.1	29	0.0003
	Distribution of observed $\gamma^2$	$\bigcirc$	0.076	22.8	17	1.4754
000	,		0.099	29.7	37	1.7943
0		$\bigcirc$	0.144	43.2	45	0.0750
80			0.064	19.2	21	0.1688
sity 600			0.167	50.1	43	1.0062
		$\bigcirc$	0.033	9.9	10	0.0010
Ο 4 .			0.160	48.0	49	0.0208
200		$\bullet$	0.108	32.4	31	0.0605
0		$\sum_{i}$	1.000	300	300	4.9716
-	0 5 10 15 20 25 30 35					
	Observed $\chi^2$					

χ2

Testing

#### Testing $\chi^2$

## Testing numbers of individuals - Computations

The population *P* splits into *b* classes with proportions  $p_i$ , so that  $\sum_{i=1}^{D} p_i = 1$ 

We take one sample *S* of size *N* out of this population How does the distribution among classes differ from the expected one?

- Expected:  $e_i = Np_i$
- *b* values  $\rightarrow 1$  indicative value:  $\chi^2 = \sum_{i=1}^{b} \frac{(s_i e_i)^2}{e_i}$



*N.B.*  $\nu = b - 1$ ; the  $\chi^2$  law does NOT depend on N!

B.B. (MMA / SBR)

### Example

In an area, 5 *protostomia* species are known to represent (in %): Pseudonereis variegata Octomeris angulosa Tetraclita serrata Patella granularis Perna perna 56 25 On one beach in this area, we count: Pseudonereis variegata Octomeris angulosa Tetraclita serrata Patella granularis Perna perna 298 149 61 32 29 Is this unexpected? **cat(**" $\ln \chi^2$  test: conformity to a known distribution $\ln$ ") pr <- c(0.56, 0.25, 0.09, 0.07, 0.03) # sum = 1.00spcnt <- read.csv("SpeciesCounting.csv".header=T)</pre> attach(spcnt) **cat**("HO: the sample comes from the population\n") # Compute expected numbers expnb = pr \* sum(sample1) # Compute  $\chi^2$  and p-value cX2 = 0for ( i in 1:length(pr) ) { cX2 = cX2 + (sample1[i]-expnb[i]) \*\* 2 / expnb[i]}  $pv = 1 - pchisq(cX2, df=length(pr)-1) \# We want P(\chi^2 > c\chi^2)$ **cat**("  $\chi^2$  =",cX2,"\np-value =",pv,"\n")

Testing  $\chi^2$ 

### # Less fun chisq.test(sample1,p=pr)

## $\chi^2$ to test for independance

We want to compare three countries for the repartition of algae<sup>7</sup>

Data:			
Country	Green	Red	Brown
France	121	403	183
Great Britain	109	347	183
Spain	94	365	155

We have no proportions from a reference population  $\rightarrow$  How can we compute

expected counts?

Hypothesis  $H_0$ : Algal group and Countries are independent factors.

$$\Rightarrow P(C_i \& G_j) = P(C_i) \times P(G_j) = \frac{nC_i}{N_{tot} \times \frac{nG_j}{N_{tot}}}; E \times p_{ij} = N_{tot} \times P(C_i \& G_j)$$

χ2

Testing

Marginal sums and expected counts				
C.	Green	Red	Brown	Sum
Fr	116.87	402.20	187.93	707
GB	105.63	363.51	169.86	639
Sp	101.50	349.29	163.21	614
Sum	324	1115	521	1960
<i>N.B.</i> $\nu = (l-1)(c-1) = 4df$				

$$\chi^{2} = \sum_{i=1}^{l} \sum_{j=1}^{c} \frac{(obs_{ij} - exp_{ij})^{2}}{exp_{ij}}$$
spfr <- c(121,403,183)
spgb <- c(109,347,183)
spes <- c(94,365,155)

chisq.test(macroa) # matrix, no p=...

<sup>7</sup> Data from T. Burel, M. Le Duff, E. Ar Gall (2019) Cah. Natur. Obs. Mar. VII:1, 1–38.

B.B. (MMA / SBR)

Stats with R

## $\chi^2$ or not $\chi^2$ : GO-term enrichment

Among 30 genes, of which 8 belong to GO-term G1, we find 10 significantly up-regulated, of which 6 belong to G1

Is the UP pool enriched in GO-term G1?

Contingency table approach:

	UP	EQ	Σc
G1	6	2	8
Gx	4	18	22
Σl	10	20	30

gomat <- matrix(c(6,4,2,18),nrow=2)
rownames(gomat) <- c("G1","Gx")
colnames(gomat) <- c("UP","EQ")
print(gomat)</pre>

chisq.test(gomat)

We get the answer: reject  $H_0$  with p = 0.013, but:

- There is a warning because of low values
- χ<sup>2</sup> says "H<sub>0</sub> can be rejected" but not "enriched / depleted" (We can look up for 1 GO-term, but extension to real-size analysis?)

G1

G1

G1 G1

G1

## Low-count-friendly approach: Fisher exact test

A very simple idea, indeed:

- enumerate all possibilities: the number of G1 among 8 can be 0, 1, ..., 8
- of for each, compute the probability: how many ways to get this pattern? under H<sub>0</sub> = UP/EQ is independent from G1/Gx Example for x = 6:
  - ► Split 30 genes into 10 + 20 (all possible patterns): Combinations of 10 genes out of 30:  $c_0 = \frac{30!}{10! \times 20!} = 30045015$
  - Split 8 genes into 6 + 2: Combinations of 6 genes out of 8:  $c_1 = \frac{8!}{6|x|^2} = 28$
  - Split 22 genes into 4 + 18: Combinations of 4 genes out of 22:  $c_2 = \frac{22!}{4|x|8|} = 7315$

The probability for this pattern is:

$$p_6 = \frac{c_1 \times c_2}{c_0} = \frac{28 \times 7315}{30045015} \approx 6.817 \times 10^{-3}$$

Sum-um all the proba of cases displaying this enrichment or more At least as much enriched: p = Σp<sub>x≥6</sub>

$$p = \Sigma p_{x \ge 6} = 7.23 \times 10^{-3}$$

### Code for the GO-term enrichment test

```
gomat <- matrix(c(6,4,2,18),nrow=2)
rownames(gomat) <- c("G1","Gx")</pre>
colnames(gomat) <- c("UP","E0")</pre>
print(gomat)
m = sum(gomat["G1",])
n = sum(gomat["Gx",])
k = sum(gomat[."UP"])
c0 = choose((m+n),k) \# (m+n)!/(n!m!)
x \leq seq(0,m)
c1 = choose(m.x)  # m!/(x!(m-x)!)
c2 = choose(n, (k-x)) # n!/((k-x)!(n-(k-x))!)
p = c1 * c2 / c0
cat("Probabilities (hypergeometric law):\n")
print(p)
cat("Sum =",sum(p),"\n")
cat("P(x>6) =".sum(p[x>=6])."\n")
```

### Lazy version:

```
dh <- dhyper(x,m,n,k)
print(dh)
cat("Sum =",sum(dh),"\n")
cat("P(x≥6) =",sum(dh[x>=6]),"\n")
```

### Super-lazy version:

```
fisher.test(gomat,alternative="greater")
```

Testing

linear dependance

### Two variables measured on the same individuals



Properties of covariance and correlation coefficient

Covariance: 
$$s_{A,B} = \frac{1}{n-1} \sum_{i=1}^{N} (a-\bar{a})(b-\bar{b})$$

An index shared by all variables, whatever their s.d. (Pearson):

$$f_{A,B} = \frac{s_{A,B}}{s_A s_B}$$



### Properties of Pearson's correlation coefficient

$$s_{A,B} = \frac{1}{n-1} \sum_{i=1}^{N} (a - \bar{a})(b - \bar{b})$$
  $r_{A,B} = \frac{s_{A,B}}{s_A s_B}$ 

Under  $H_0$  = "Two uncorrelated variables":

$$r\sqrt{rac{n-2}{1-r^2}} \mapsto t_{\nu=n-2}$$

if variables are normal or sample is large.

### Example

Correlation between location of algae, and their content in Fucose + sulfate<sup>8</sup>

Species	Emersion(%)	Fucose (% D.W.)
Ascophyllum nodosum	9.5	16.1
Bifurcaria bifurcata	26.8	12.6
Fucus cerranoides	63.7	22.6
Fucus serratus	30.1	10.9
Fucus spiralis	81.6	21.4
Fucus vesiculosus	49.6	17.2
Laminaria digitata	12.8	3.6
Pelvetia canaliculata	98.5	34.9

# Data

fucdat <- read.table("Fucose.tab",header=T)
attach(fucdat)</pre>

```
# Parameters
sde = sd(Emersion) ; sdf = sd(Fucose)
co = cov(Emersion,Fucose)
r = co/(sde*sdf)
cat("Correlation r =",r,"\n")
```

```
# Test
free = nrow(fucdat)-2
t = r * sqrt(free/(1-r*r))
```



# less fun
r = cor(Emersion,Fucose)
cat("Correlation r =",r,"\n")
cor.test(Fucose,Emersion,alternative="g")

#### # Plot

plot(Emersion,Fucose, main="Fucose vs location", xlab="Emersion duration (%)", ylab="Fucose + Sulfate (%D.W.)" )

cat("p-value for positive correlation =",dt(t,df=free),"\n")

<sup>8</sup> Data from B. Kloareg (1991), Bull. Soc. Bot. Fr. 138:3-4, 305-318.

### Linear regression

Assuming dependance between variables makes sense...

Linear correlation means  $y = ax + b + \epsilon$  where  $\epsilon =$  "residual" Estimates:  $\hat{a}$  and  $\hat{b}$  for which  $\hat{y} = \hat{a}x + \hat{b}$  predicts y with the lowest error.  $\Rightarrow$  Define error criterion:  $SSR = \sum_{n} \epsilon^2 = \sum_{n} (y - \hat{y})^2$ 



Good news: no need to try all possible values!

$$\hat{a} = r \frac{s_y}{s_x}$$
;  $\hat{b} = \bar{y} - \hat{a}\bar{x}$ 

Confidence intervals:

$$s_{a} = s_{SSR} / \sqrt{\sum (x_{i} - \bar{x})^{2}} \qquad \Delta \hat{a} \mapsto s_{a} t_{\nu = n - 2}$$
  
$$s_{b} = s_{SSR} \sqrt{\sum x_{i}^{2} / n \sum (x_{i} - \bar{x})^{2}} \qquad \Delta \hat{b} \mapsto s_{b} t_{\nu = n - 2}$$



#### linear dependance

## Example

How do tides impact phytoplankton biomass in Penze?<sup>9</sup>

 $\rightarrow$  Correlate chorophyl content with water level.

```
# Check linear correlation
rc = cov(H,Ch1)/(sd(H)*sd(Ch1))
```

```
# Estimates for a and b
ac = rc * sd(Chl) / sd(H)
bc = mean(Chl) - ac * mean(H)
cat("Chl = ",ac," H +",bc,"\n")
```

```
# Use the linear model function
datac <- data.frame(chloro)
regc <- lm(Chl~H,datac)</pre>
```

```
cat("Regression coefficients:\n")
print(regc$coef)
cat("Confidence interval:\n")
confint(regc,level=0.95)
```

```
# Plot
minh=1.5
maxh=5.0
minc=1.0
maxc=3.0
```



```
sprintf("a=%.3f b=%.3f\nSSR = %.4f",ac,bc,SSRc))
```

```
cat(" Regression residuals:\n")
print(regc$residuals)
SSRc = sum(regc$residuals**2)
cat(" SSR =",SSRc,"\n")
```

<sup>9</sup> Data from C. Riaux and J.L. Douvillé (1980), Estuar. Coast. Mar. Sci 10, 85-92.

B.B. (MMA / SBR)

### Correlation and prediction

Regression = modeling the relation between *x* and *y* 

$$\Rightarrow$$
 knowing *x*, predict *y* as  $\hat{y} = \hat{a}x + \hat{b}$ 

L



BUT uncertainty  $\Delta \hat{a}$  and  $\Delta \hat{b} \Rightarrow \Delta \hat{y}$ 

$$\Delta \hat{y} \quad \mapsto \quad \left(s_{SSR} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum (x_i - \bar{x})^2}}\right) t_{\nu = n - 2}$$

N.B. the more *x* differs from  $\bar{x}$ , the highest  $\Delta \hat{y}$ 



#### linear dependance

### Predict Y vs predict X

Regression = modeling the relation between *x* and *y* as y = ax + b

 $\Rightarrow \text{ knowing } x, \text{ predict } y \text{ for } x \text{ within a reasonable range.}$  $\Rightarrow \text{ knowing } y, \text{ predict } x?$ 



B.B. (MMA / SBR)

### More than one factor



Question: does the value of *X* depend on...

- factor A?
- factor B?
- factor C?
- an interaction between A and B?
- an interaction between A and C?

• ...?

Model: for each individual i,

 $X_i = aA_i + bB_i + cC_i + X_0 + \epsilon_i$ 

Problem 1: knowing  $A_i$ ,  $B_i$ ,  $C_i$  for all i, estimate a, b, c,  $X_0$  so that the  $X_i$  are well predicted by:

$$\hat{X}_i = \hat{a}A_i + \hat{b}B_i + \hat{c}C_i + \hat{X}_0$$



Problem 2: Which factor(s) display(s) a significant influence?

### Simple answers

Answer to problem 1 = find coefficients:

• Simply use linear regression: 1m can cope with multiple factors

```
multimodel <- lm(X \sim A + B + C)
```

Answer to problem 2 = effect of each factor:

• Simply use AnOVa: anova can cope with multiple factors anova(multimodel)

Bonus = combined effect of factors:

• Simply replace + by \*

```
combomodel <- lm( X ~ A * B * C )
anova(combomodel)</pre>
```

## A real, thus less simple, situation

Factors having an effect on brown alga fecundity<sup>10</sup>.

Question: how does fecundity depend on the three following factors:

- species considered = Laminaria ochroleuca, Sacchoriza polyschides ;
- water temperature = 10°C, 15°C, 25°C ;
- time = 0, 2, ... , 26 days

```
\rightarrow How about something like: AlFec <- Im(Fec \sim Spe * Tmp * Day)
anova(AlFec)
```

Problem: Temperature and Time are quantitative, Species are qualitatitive (categories) Solution: Turn all variables to qualitative, *i.e.* factors!

```
fectbl <- read.table("SpeDayTmpFec.tsv",header=T)
fectbl[,1:3] <- lapply(fectbl[,1:3],as.factor)
model <- aov(Fec ~ Spe * Day * Tmp, data = fectbl)
summary(model)</pre>
```



Bonus: because model is made by aov – not lm – *post-hoc* test applies:

```
TukeyHSD(model,c("Spe","Day","Tmp"))
```

<sup>&</sup>lt;sup>10</sup>data very approximately reconstructed from T.R. Pereira, A.H. Engelen, G.A. Pearson, E.A. Serrão, C. Destombe, M. Valero (2011) *Cah. Biol. Mar.* **52**:395-403.

## Non-linear parameter optimization

Example: Temperature and Bacterial Production<sup>11</sup> Bacterial production follows a non-linear law:



$$BP(T) = BP_{max} \left( \frac{T_{max} - T}{T_{max} - T_{opt}} \right)^{\beta} \exp\left( -\beta \left( \frac{T_{max} - T}{T_{max} - T_{opt}} - 1 \right) \right)$$

with  $BP_{max}$ ,  $T_{max}$ ,  $T_{opt}$  and  $\beta$ : parameters to estimate.

```
# create the function
metatemp <- function(temp,BPmax,Tmax,Topt,beta) {</pre>
  BPmax*((Tmax-temp)/(Tmax-Topt))**beta*
   exp(-beta*(((Tmax-temp)/(Tmax-Topt))-1))
 read data and build the data frame
Ħ
tbptbl <- read.table("TempBPresp.tab".header=T)</pre>
attach(tbptbl)
databp <- data.frame(tbptbl)</pre>
# optimize parameters
opttbp <- nls(
    BP~metatemp(Temp.BPmax.Tmax.Topt.beta).
    data=databp,
    start=list(BPmax=max(BP),Tmax=max(Temp),
                Topt=27.beta=1) )
opc <- coef(opttbp)</pre>
print(opc)
```



# plot

tpx < -seq(2, 34, length.out = 100)bpy<-metatemp(tpx,opc[1],opc[2],opc[3],opc[4])</pre> plot(tpx,bpy,type="1",lty=2,lwd=3,col="red", main="BP - temp. response", xlab="Temperature (°C)",ylab="BP (mqC/m<sup>2</sup>/h)") points(databp,pch=16,col="darkgreen",cex=2)

11 Data from C. Hubas, L.F. Artigas, D. Davoult (2007). Mar Ecol Prog Ser 344:39-48.

B.B. (MMA / SBR)
## Workflow of statistic test

