

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

Seminar III: R/Bioconductor

Misc. Stats and affy Quality Control.

José Víctor Moreno Mayar
jmoreno@lcg.unam.mx

LCG - UNAM

August - December, 2009

Misc. Stats and affy Control Quality.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

① Packages for this class.

② lowess.

③ loess.

④ Affymetrix Chips.

⑤ AffyBatch

⑥ phenoData

Misc. Stats and affy Control Quality.

7 Affy quality Control.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

Packages for this class.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

For this class we will use two functions already included in R, and the affy package.

```
> source("http://bioconductor.org/biocLite.R")  
> biocLite("affy")
```

Some theory.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- It is a method for making smooth regressions of scatter plots.
- It uses the **least squares** method.
- Let us recall the method.

Least Squares.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- It consists, mainly, on fitting a curve to a given set of points.
- The total sum of the residuals has to be minimized.¹
- The function depends on the degree of the polynomial that should be fitted.
- A system of linear equations is gotten so the coefficients of the regression curve are calculated.

¹The square of the error is used, so it is called least squares

Locally Weighted Regression.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The difference here is the weighting function $W(x)$ used for each x_i .
- So each pair (x_i, y_i) has a different effect on the regression.
- Even more, $W(x)$ is calculated based on the neighborhood of x_i .

The Weighting Function.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- An initial set of weights $w_k(x_i)$ is calculated for each x_i .
- An initial fitted value \hat{y}_i is calculated for each x_i .
- This is made by means of **weighted least squares**.
- The point is fitted to a d th degree polynomial function (usually, $d = 1$).

f and The Smoothness Assumption.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The fact that the neighbors of an x_i may be used in calculating its set of weights, makes an assumption of smoothness.
- So, you can select how wide will the neighborhood be.
- This is made through the parameter f .
- Different values of f may be taken from the interval $(0,1]$.
- The larger the f , the wider the neighborhood and the smoother the points will be.

Robust Lowess.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The way to make this regression robust is by calculating a better \hat{y}_i .
- The regression is guarded against deviant points, which might distort the smoothed point by reiterating the procedure.
- This time, the computing of the \hat{y}_i will be made with a different set of weights δ_i .
- δ_i is calculated based on the size of the residual $y_i - \hat{y}_i$.
 - ▶ Large residuals result in small weights and small residuals result in large weights.

Lowess in R.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The function in R to generate these curves is `lowess()`.
- The arguments to be passed to the function are the two related (or not) variables, `f`, which is the span of the smoothing process and the number of times the process should be iterated.
- Time for some practice.

An example.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

Let us create a data frame with two variables x and y where:

- x are the first 151 consecutive natural numbers.
- $y_i = .2x_i + \epsilon_i$.
- ϵ_i is a random sample taken from a normal distribution with $\mu=0$ and $\sigma=1$.
- Make an xyplot of it.
- Adjust 5 curves with different values for f $f=.01, f=.2, f=.5, f=.8, f=.99$, and plot them.
- Compare it to a curve made with `lm()`.

An example.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

Creating the data frame.

```
> x <- seq(from = 0, to = 150, by = 1)
> error <- rnorm(1000, 0, 1)
> error151 <- sample(error, 151,
+     replace = T)
> y <- 0.02 * x + error151
> xy <- data.frame(x, y)
```

An example.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

Making the regressions.

```
> lo.xy.01 <- lowess(xy$x, xy$y,  
+      f = 0.01)  
> lo.xy.2 <- lowess(xy$x, xy$y, f = 0.2)  
> lo.xy.5 <- lowess(xy$x, xy$y, f = 0.5)  
> lo.xy.8 <- lowess(xy$x, xy$y, f = 0.8)  
> lo.xy.99 <- lowess(xy$x, xy$y,  
+      f = 0.99)  
> lr <- lm(y ~ x, data = xy)
```

An example. (xyplot and regressions.)

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

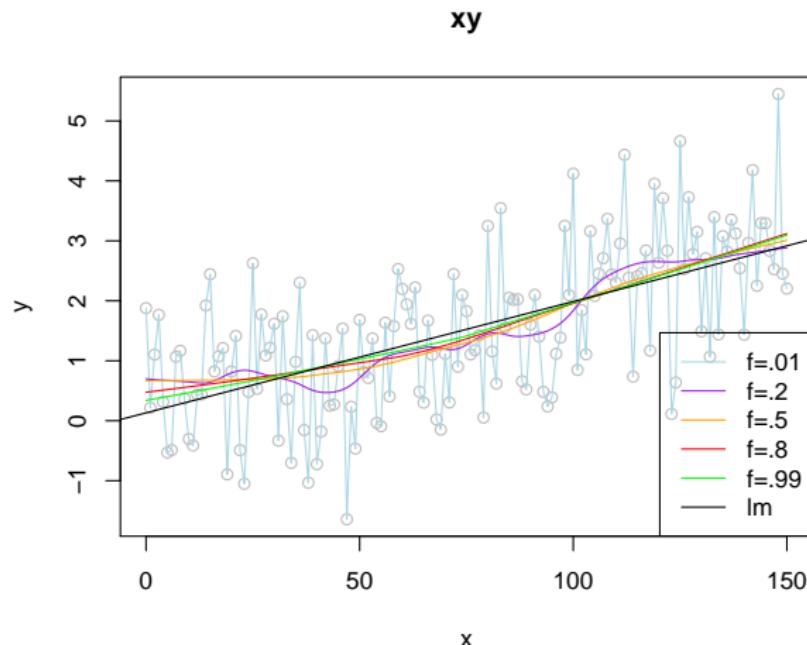
loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.



Pros and Cons.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- What does the function `lowess()` return?²
- What are its cons when compared to `lm()`?³
- Which one do you consider to be the best fit?
- Do you notice any convergence?

²Use `class()`.

³Which functions could be used on an object of class `lm`?

Choosing f .

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- f should be chosen in order for it to maximize the smoothing.
- It should not make the pattern in the data diffuse or distorted.
- Most of the times, f should be chosen in the interval (.2,.8).
- When you do not know which is the best value for f , .5 should make the job, for an initial exploratory analysis.

loess.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch
phenoData

Affy quality
Control.

- loess is a **newer** version of lowess.
- It implements the same local fitting method as lowess.
- The iterative method is optional.
- loess uses the formula notation.
- Here, you can decide the degree of the polynomial to which the data should be fitted.
- If you do not rely on the least squares approximation method, you can select some other one.
- An important argument of the function is **loess.control**.

loess.control.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- loess.control helps you tune different parameters of the lowess computation.
- surface -> this regression may be used to extrapolate values based on the data.
- statistics, trace.hat -> to set these parameters to approximate would be useful if there are several points.
- Here, you can select the number of iterations to be made, as well.

Extrapolating with loess.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Let us do an extrapolation with the surface argument of the loess.control set to "direct".

```
> xy.lo2 <- loess(y ~ x, xy, control = loess.control(  
> pred2 <- predict(xy.lo2, data.frame(xs = seq(50,  
+      200, 1)), se = TRUE)
```

Plotting the Extrapolation.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

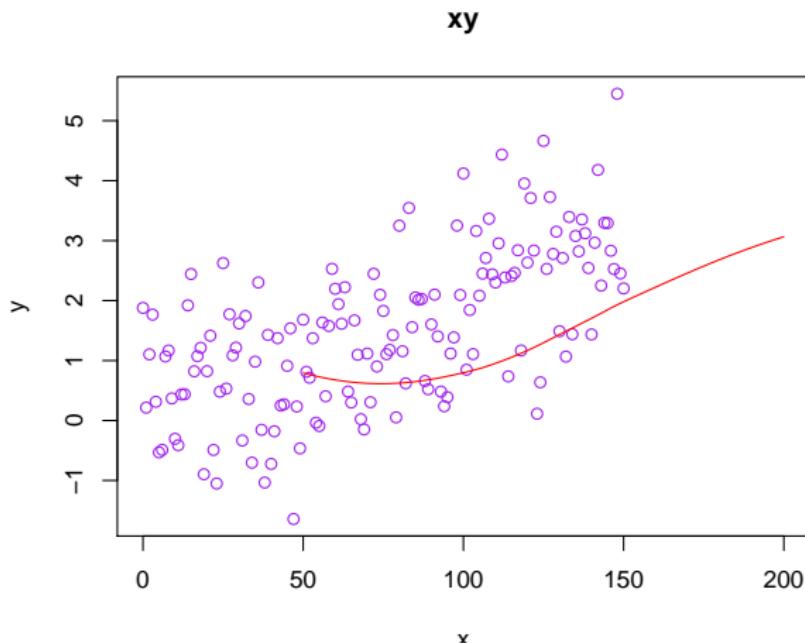
loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.



About the chips.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Let us study some important aspects of the Affymetrix oligonucleotide arrays tech.
- The goal is to probe an RNA sample (target) with different oligonucleotide probes.⁴
- Each feature is called a **probe pair**.

⁴25 bp long

About the chips.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Each probe pair contains a **perfect match (PM)** and a **mismatch (MM)** probe.
 - ▶ MM is used to measure the amount of **non-specific binding**.⁵
- A **probe pair set** is made up of all the PMs and MMs related to a common **affyID**.⁶

⁵It is created by changing the 13th base of the PM.

⁶An affyID is related to a gene or a gene fraction represented on the array.

About the chips.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- You can access to the PM, MM and probe names data with the `pm()`, `mm()`, `probeNames()` methods, for an `AffyBatch` object.⁷
- Another important function is `geneNames()`, which extracts unique affyIDs from an AffyBatch object.

⁷We will learn more about it, later.

Reading .CEL files.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Affy chips are read into .CEL files which are read into R with the `ReadAffy()` function.
- First of all, you should have all your .CEL files in your working directory. You can check or set it with the `getwd()` and `setwd()` functions.
- All you have to do, is to keep your data into an R object.

```
> library(affy)
> data <- ReadAffy()
> data
```

Reading .CEL files.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

AffyBatch object

size of arrays=732x732 features (18 kb)

cdf=HG-U133A_2 (22277 affyids)

number of samples=6

number of genes=22277

annotation=hgu133a2

notes=

AffyBatch Objects.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- What is the class of our object *data*?
- You can explore it by means of the **slotNames()** function.⁸
- If you want to access a particular slot, you can do it with the following syntax.

```
> slotNames(data)
```

```
[1] "cdfName"  
[2] "nrow"  
[3] "ncol"  
[4] "assayData"  
[5] "phenoData"  
[6] "featureData"  
[7] "experimentData"  
[8] "annotation"
```

AffyBatch Objects.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

```
[9] "protocolData"  
[10] ".__classVersion__"  
  
> slot(data, "nrow")  
  
[1] 732
```

⁸S3 and S4 objects can be accessed this way

phenoData.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Microarray expression data is useless without metadata⁹, so there must be a description of what is being assessed in each chip.
- You can find the description in the **phenoData** slot of an AffyBatch object.
- You can even assign a phenoData file to your data manually.
- This will be useful when testing for differentially expressed genes.

⁹Dr. Salt made a point on it last Monday.

Assigning a phenoData file

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Download the file pdata1422.txt and move it to your working directory.
- Then, read it with the `read.AnnotatedDataFrame()`.
- Finally, put it into your AffyBatch object.

Assigning a phenoData file.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

```
> pData(slot(data, "phenoData"))
```

	sample
AF14.CEL	1
AF15.CEL	2
AF16.CEL	3
AF6.CEL	4
AF7.CEL	5
AF8.CEL	6

```
> pd <- read.AnnotatedDataFrame(filename = "pdata1422"
+           header = T)
> pData(pd)
```

Assigning a phenoData file.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

	Source.Name
AF16.CEL	PROX1_siRNA-2_replicate2
AF7.CEL	PROX1_siRNA-1_replicate1
AF14.CEL	GFP_siRNA_replicate2
AF8.CEL	PROX1_siRNA-2_replicate1
AF15.CEL	PROX1_siRNA-1_replicate2
AF6.CEL	GFP_siRNA_replicate1

```
> slot(data, "phenoData") <- pd  
> pData(slot(data, "phenoData"))
```

Assigning a phenoData file.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

	Source.Name
AF16.CEL	PROX1_siRNA-2_replicate2
AF7.CEL	PROX1_siRNA-1_replicate1
AF14.CEL	GFP_siRNA_replicate2
AF8.CEL	PROX1_siRNA-2_replicate1
AF15.CEL	PROX1_siRNA-1_replicate2
AF6.CEL	GFP_siRNA_replicate1

Getting Microarray Data.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- As you know, you can use the **ArrayExpress** package to get expression data.
- I got my .CEL files as follows:

```
> library(ArrayExpress)
> Data = ArrayExpress("E-MEXP-1422",
+           save = T)
```

- This AffyBatch object already has its phenoData complete, so, if you do not know much about the sample, ArrayExpress is a must.

image.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- Sometimes, the chip will be damaged or some dust will fall on it, so some expression values obtained from it, will not be reliable.
- You can see the chips, as they were scanned, to check for this.

image.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

```
> par(mfrow = c(2, 3))  
> image(Data)  
> par(mfrow = c(1, 1))
```

hist.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The **hist** function is well suited to deal with AffyBatch objects, as well.
- Make a histogram in order check if the intensities are distributed similarly among the arrays.
- You can infer the need for **normalization** between arrays with this plot.¹⁰

¹⁰We will discuss this in the next class.

hist.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

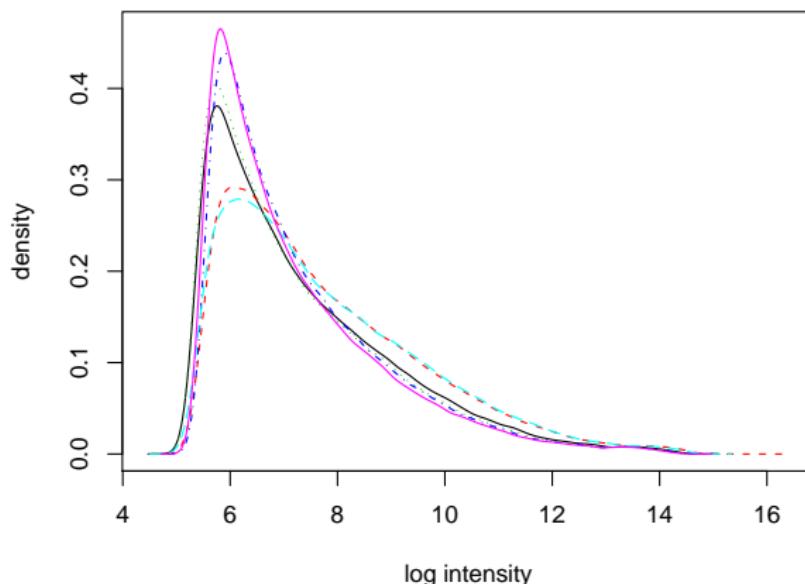
loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.



M vs A plot.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The log ratio of the intensities is plotted on the vertical axis.
- The average of the log intensities is plotted on the horizontal axis.
- This plots offer a way of making pairwise graphical comparison of intensity data.
- Problems in replicate sets of arrays may be assessed with these graphs.
- We will use our Data object as it is ordered according to the replicates.¹¹

¹¹See the phenoData.

M vs A plot.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

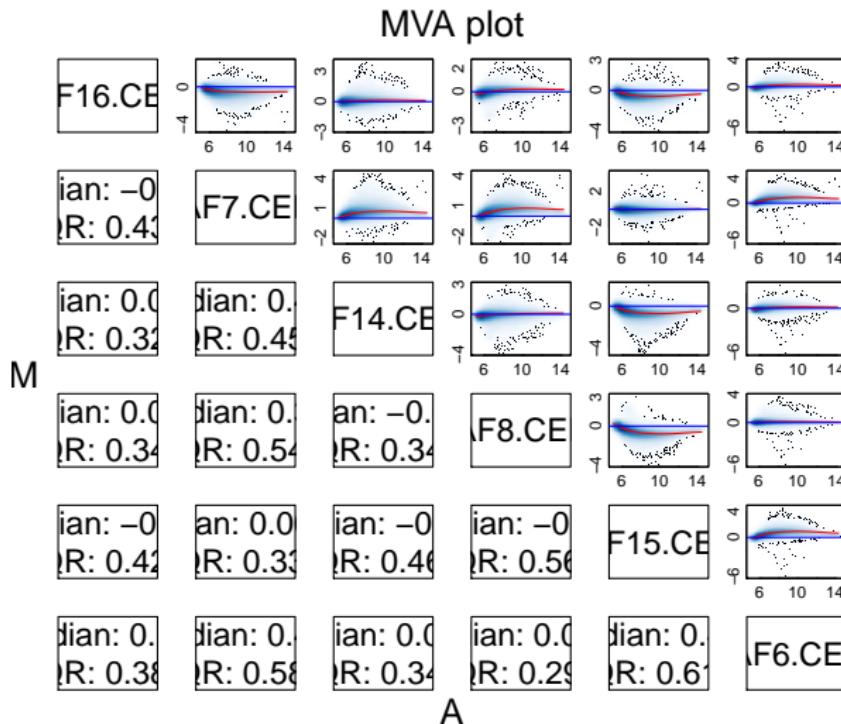
loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.



M vs A plot.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- The argument pairs sets the function to make pairwise comparisons.
- What would happen with pairs=F?¹²
- What can you notice about the red regression and the deviant points?¹³

¹²The chip is compared to a reference chip.

¹³We need to normalize, wait for the next class.

RNA Degradation.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

- As you already know, different fragments of the same RNA are represented in different features of the array.
- An artefact of the microarray technique is that features representing sites closer to the 5' end of the RNA are less intense.
- This is because RNA degradation in the assay **starts at the 5' end** of the molecule.
- There is a way to check for the RNA degradation rate in the assay.

RNA Degradation.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

```
> degradation <- AffyRNAdeg(Data)
> names(degradation)

[1] "N"                  "sample.names"
[3] "means.by.number" "ses"
[5] "slope"              "pvalue"
```

- There is a way to make a summary of this object.¹⁴
- Now, we will plot **degradation curves** for this data.

¹⁴Use the `summarizeAffyRNAdeg()` function.

RNA Degradation Plots.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

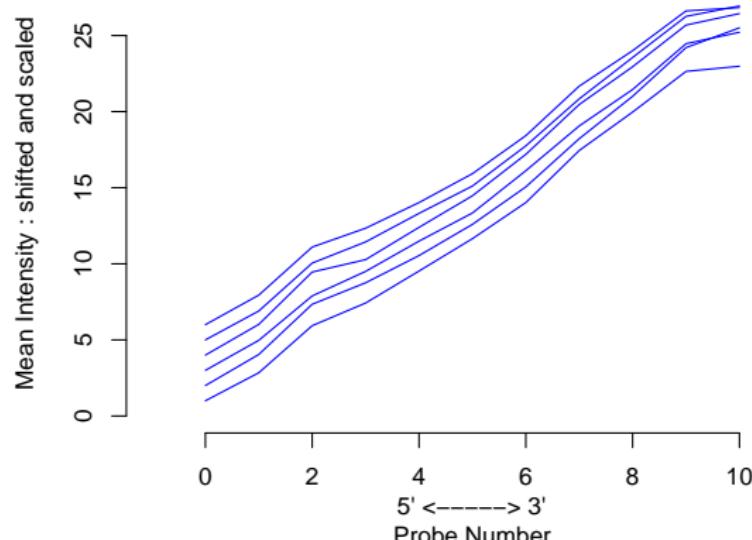
Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

RNA degradation plot



Pending Points.

Seminar III:
R/Bioconductor

Víctor Moreno

Class
Overview.

Packages for
this class.

lowess.

loess.

Affymetrix
Chips.

AffyBatch

phenoData

Affy quality
Control.

Next class, we will learn about microarray normalization,
finding DEGs and the SpeCond package.