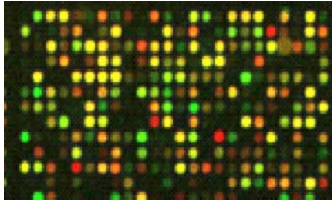


## Statistics for cDNA Microarrays

### Image Analysis: Quality assessment and exploratory data analysis

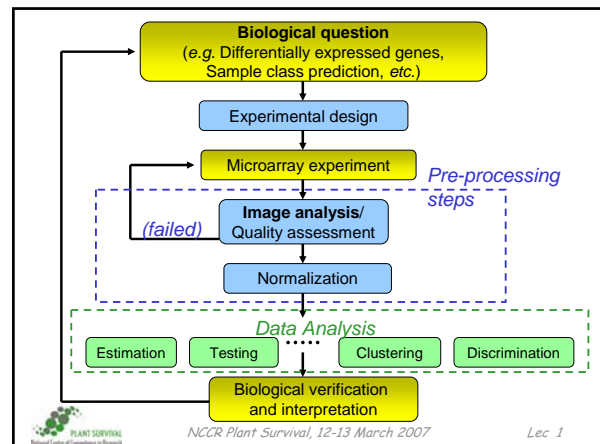


<http://www.isrec.isb-sib.ch/~darlene/NCCR-PS/>



NCCR Plant Survival, 12-13 March 2007

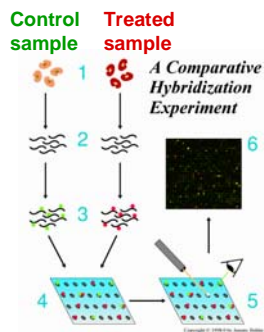
Lec 1



NCCR Plant Survival, 12-13 March 2007

Lec 1

## Microarray Experiment



NCCR Plant Survival, 12-13 March 2007

Lec 1

## Measures of center: Mean

- The *mean* value of a variable is obtained by computing the total of the values divided by the number of values
- Appropriate for data patterns that are fairly *symmetrical*
- It is sensitive to presence of outliers, since all values contribute equally



NCCR Plant Survival, 12-13 March 2007

Lec 1

## Measures of center: Median

- The *median* value of a variable is the number having 50% (half) of the values smaller than it (and the other half bigger)
- It is NOT sensitive to presence of outliers, since it 'ignores' almost all of the data values
- The median is thus usually a more appropriate summary for *skewed* (asymmetrical) data patterns



NCCR Plant Survival, 12-13 March 2007

Lec 1

## Quantification of Expression

For each spot on the slide, calculate

$$\text{Red intensity} = R_{fg} - R_{bg}$$

(fg = foreground, bg = background) and

$$\text{Green intensity} = G_{fg} - G_{bg}$$

and combine them in the log (base 2) ratio

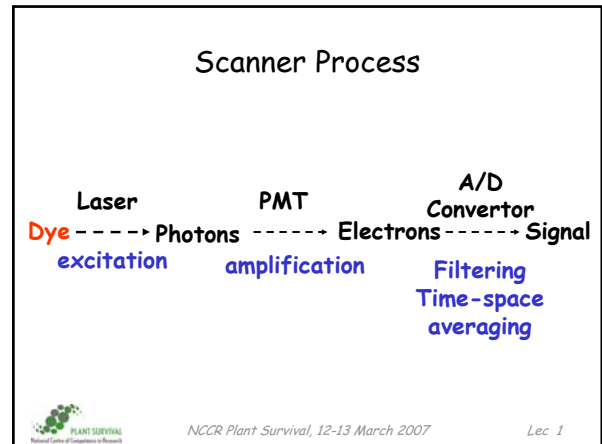
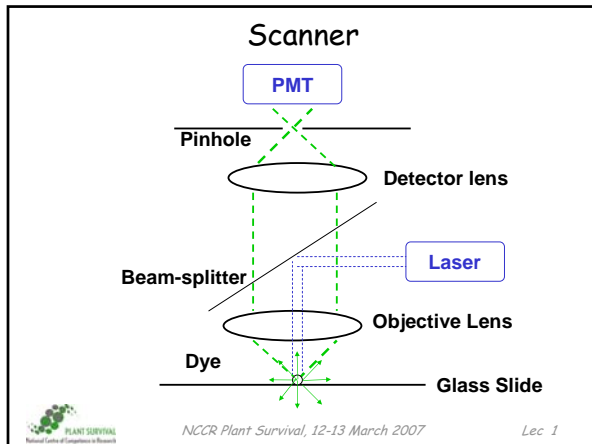
$$\text{Log}_2(\text{Red/Green})$$

Often, fg = mean and bg = median of relevant pixel intensities



NCCR Plant Survival, 12-13 March 2007

Lec 1



### Images from Scanner

- Resolution
  - standard 10µm [currently, best ~ 5µm]
  - 100µm spot on chip = 10 pixels in diameter
- Image format
  - TIFF (tagged image file format) 16 bit (65,536 levels of gray) - also other formats
  - 1cm x 1cm image at 16 bit = 2Mb
- Separate image for each fluorescent sample
  - channel 1, channel 2, etc.

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Images : examples

Pseudo-color overlay

Spot color	Signal strength	Gene expression
yellow	Control = Treated	unchanged
red	Control < Treated	induced
green	Control > Treated	repressed

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Steps in Images Processing

- *Addressing* (or *Gridding*)
  - Assigning coordinates to each spot
- *Segmentation*
  - *Classification of pixels* as either foreground (signal) or background
- *Information Extraction*
  - Foreground fluorescence intensity pairs (*R,G*)
  - Background intensities
  - Quality measures

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Addressing

This is the process of assigning coordinates to each of the spots.

Automating this part of the procedure permits high throughput analysis.

4 by 4 grids  
19 by 21 spots per grid

NCCR Plant Survival, 12-13 March 2007 Lec 1

## Addressing — Registration

NCCR Plant Survival, 12-13 March 2007
Lec 1

## Addressing (I)

- Basic structure of images **known** (determined by the arrayer)
- Parameters to address spot positions
  - Separation between rows and columns of grids
  - Individual translation of grids
  - Separation between rows and columns of spots within each grid
  - Small individual translation of spots
  - Overall position of the array in the image

NCCR Plant Survival, 12-13 March 2007
Lec 1

## Steps in Images Processing

- *Addressing (or Gridding)*
  - Assigning coordinates to each spot
- *Segmentation*
  - *Classification of pixels* as either foreground (signal) or background
- *Information Extraction*
  - Foreground fluorescence intensity pairs ( $R, G$ )
  - Background intensities
  - Quality measures

NCCR Plant Survival, 12-13 March 2007
Lec 1

## Segmentation

- Classification of pixels as foreground or background
  - fluorescence intensities are calculated for each spot as measure of transcript abundance
- Production of a *spot mask*: set of foreground pixels for each spot

NCCR Plant Survival, 12-13 March 2007
Lec 1

## Segmentation methods in some programs

Fixed circle	ScanAlyze, GenePix, QuantArray
Adaptive circle	GenePix, Dapple, SignalViewer (uses ellipse)
Adaptive shape	Spot, region growing and watershed
Histogram	ImaGene, QuantArray, DeArray and adaptive thresholding

NCCR Plant Survival, 12-13 March 2007
Lec 1

## Fixed circle segmentation

- Fits a circle with a *constant diameter* to all spots in the image
- Easy to implement
- The spots should be of the same shape and size

*May not be good for this example*

NCCR Plant Survival, 12-13 March 2007
Lec 1

### Adaptive circle segmentation

- The circle diameter is estimated *separately* for each spot

*Dapple* finds spots by detecting edges of spots (second derivative)

- Problematic* if spot exhibits oval shapes

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Limitation of circular segmentation

- Small spot
- Not circular

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Limitation of fixed circles

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Adaptive shape segmentation

- Specification of *starting points* or *seeds*
  - Bonus: already know geometry of array
- Regions grow outwards from the seed points preferentially according to the difference between a pixel's value and the running mean of values in an adjoining region

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Histogram segmentation

- Choose target mask *larger than any spot*
- Fg and bg intensities determined from the histogram of pixel values for pixels within the masked area
- Example : QuantArray
  - Background : mean between 5th and 20th percentile
  - Foreground : mean between 80th and 95th percentile
- May not work well when a large target mask is set to compensate for variation in spot size

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1

### Steps in Images Processing

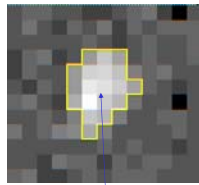
- Addressing* (or *Gridding*)
  - Assigning coordinates to each spot
- Segmentation*
  - Classification of pixels as either foreground (signal) or background
- Information Extraction*
  - Foreground fluorescence intensity pairs ( $R, G$ )
  - Background intensities
  - Quality measures

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1

## Information Extraction

- *Spot Intensities*
  - mean of pixel intensities
  - median of pixel intensities
  - Pixel variation (e.g. IQR)
- *Background values*
  - None
  - Local
  - Constant (global)
  - Morphological opening
- *Quality Information*



Take the average



NCCR Plant Survival, 12-13 March 2007

Lec 1

## Spot 'foreground' intensity

- The total amount of hybridization for a spot is proportional to the *total fluorescence* generated by the spot
- Spot intensity = sum of pixel intensities within the spot mask
- Since later calculations are based on *ratios* between Cy5 and Cy3, we compute the average\* pixel value over the spot mask
  - \* *alternative*: ratios of medians may be better than means if bright specks present



NCCR Plant Survival, 12-13 March 2007

Lec 1

## Background intensity

- The measured fluorescence intensity includes a contribution of *non-specific hybridization* and *other chemicals* on the glass
- Fluorescence from regions not occupied by DNA should be *different* from regions occupied by DNA
  - one solution is to use *local negative controls* (spotted DNA that should not hybridize)



NCCR Plant Survival, 12-13 March 2007

Lec 1

## BG: None

- Do not consider the background
  - Probably not accurate in many cases, but may be better than some forms of local background determination



NCCR Plant Survival, 12-13 March 2007

Lec 1

## BG: Local

- Focus on small regions surrounding the spot mask
- Median of pixel values in this region
- Most software implements such an approach



Scanalyze



ImaGene



Spot, GenePix

- By ignoring pixels immediately surrounding the spots, bg estimate is *less sensitive* to the performance of the segmentation procedure



NCCR Plant Survival, 12-13 March 2007

Lec 1

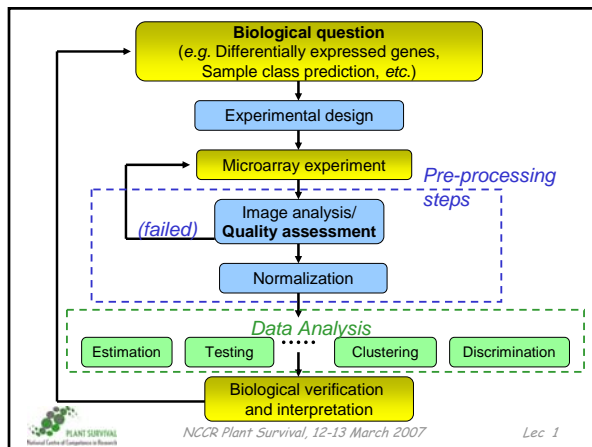
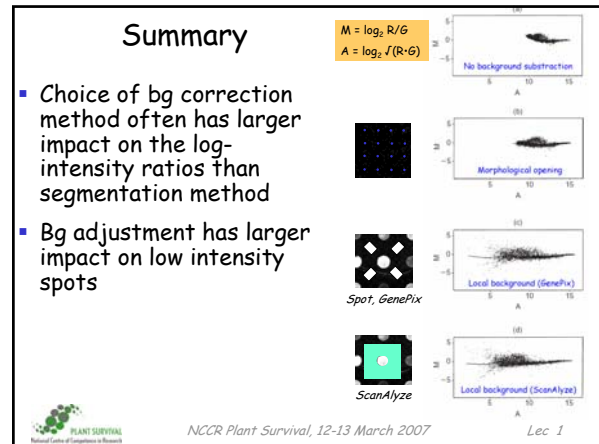
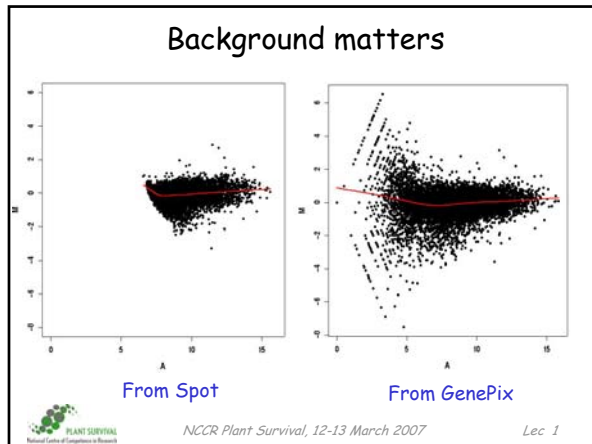
## BG: Constant

- *Global method* which subtracts a constant background for all spots
- Some evidence that the binding of fluorescent dyes to 'negative control spots' is lower than the binding to the glass slide
- → More meaningful to estimate background based on a *set of negative control spots*
  - If no negative control spots: approximation of the average background = third percentile of all the spot foreground values



NCCR Plant Survival, 12-13 March 2007

Lec 1



### Preprocessing: Data Visualization, Exploratory Data Analysis (EDA)

- Was the experiment a success?
- Are there any specific problems?
- What analysis tools should be used?

NCCR Plant Survival, 12-13 March 2007      Lec 1

### Microarray data preprocessing

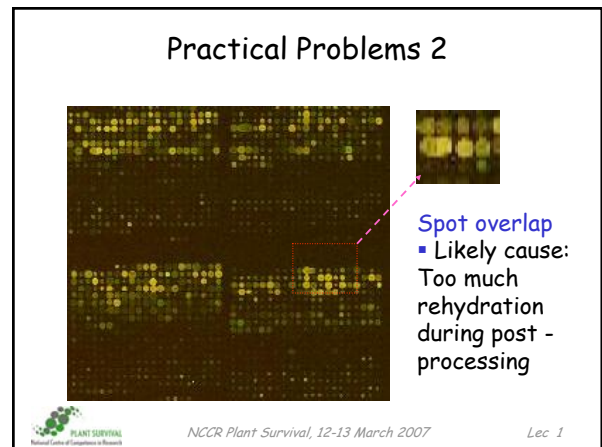
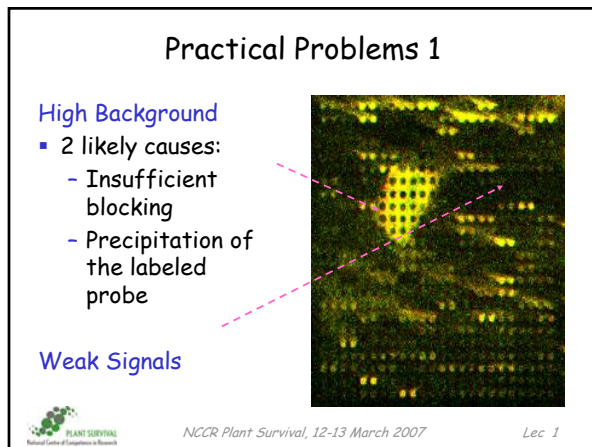
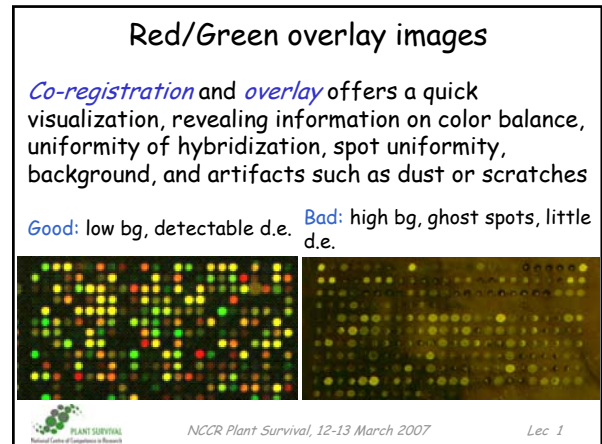
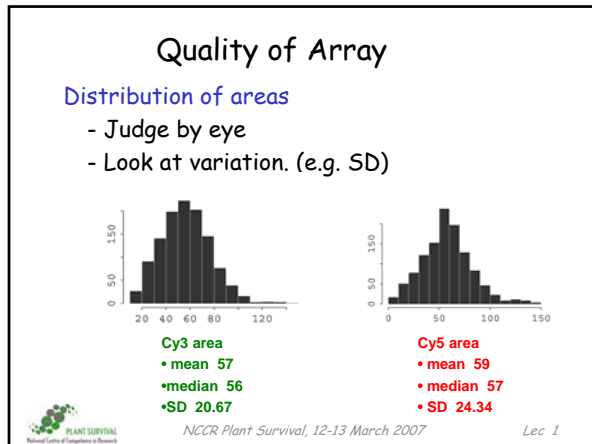
- Primary (dual channel) microarray data are pixel-level values from the *two tiff files*
- Prior to analysis, usually microarray data are *preprocessed* so that there is a *single value* for each spot
- The main preprocessing steps are
  - Image analysis (we just looked at this)
  - Normalization (to come after the break)
- We assume now that image analysis has been carried out, so that pixel values are summarized as *Rf, Rb, Gf, Gb*

NCCR Plant Survival, 12-13 March 2007      Lec 1

### Quality Measurements

- Array**
  - Correlation between spot intensities
  - Percentage of spots with no signals
  - Distribution of spot signal area
- Spot**
  - Signal / Noise ratio
  - Variation in pixel intensities
  - Identification of "bad spot" (spots with no signal)
- Ratio (2 spots combined)**
  - Circularity
- Flag or weight** spots based on these (or other) criteria

NCCR Plant Survival, 12-13 March 2007      Lec 1



### Artifacts in microarrays

- We are interested in finding true *biologically meaningful differences* between sample types
- Due to other sources of systematic variation, there are also usually *artifactual differences*
- Sources of artifacts include:
  - print tips - differences in subarrays
  - plate effects - differences in rows within subarray
  - batch effects
  - hybridization artifacts

PLANT SURVIVAL  
National Centre of Excellence in Research

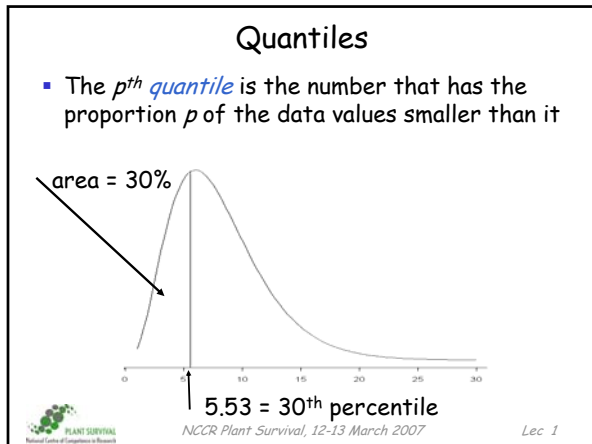
NCCR Plant Survival, 12-13 March 2007 Lec 1

### Looking for artifacts

- Exploratory data analysis (EDA) is an important component of microarray data preprocessing
- EDA involves identifying data artifacts
- We will use several types of plots for data visualization, primarily
  - *image plots*
  - *scatterplots*
  - *boxplots*
  - *spatial plots*

PLANT SURVIVAL  
National Centre of Excellence in Research

NCCR Plant Survival, 12-13 March 2007 Lec 1



- ### Five-number summary and boxplot
- The 25<sup>th</sup> ( $Q_1$ ), 50<sup>th</sup> (median), and 75<sup>th</sup> ( $Q_3$ ) percentiles divide the data into 4 equal parts; these special percentiles are called *quantiles*
  - An overall summary of the distribution of variable values is given by the five values:
 

Min,  $Q_1$ , Median,  $Q_3$ , and Max
  - A *boxplot* provides a visual summary of this five-number summary
- NCCR Plant Survival, 12-13 March 2007

