## A New Contamination Model

### Ruben Zamar Department of Statistics, UBC

October 21, 2015

# A NEW ROBUSTNESS MODEL

-

3

## More General (and Realistic) Robustness Model

• "Classical" robust methods work as follows

- "Classical" robust methods work as follows
  - Identify unusual data points in the dataset (rows in the data table)

- "Classical" robust methods work as follows
  - Identify unusual data points in the dataset (rows in the data table)
  - Downweight the unusual data cases

- "Classical" robust methods work as follows
  - Identify unusual data points in the dataset (rows in the data table)
  - Downweight the unusual data cases
- Important assumption underlying classical robust procedures

- "Classical" robust methods work as follows
  - Identify unusual data points in the dataset (rows in the data table)
  - Downweight the unusual data cases
- Important assumption underlying classical robust procedures
  - Percentage of unusual data points is relatively small

- "Classical" robust methods work as follows
  - Identify unusual data points in the dataset (rows in the data table)
  - Downweight the unusual data cases
- Important assumption underlying classical robust procedures
  - Percentage of unusual data points is relatively small
  - Hopefully way below 50%

• Typical data table (before)

D Α Т Α Т Α В L Ε

#### • In some applications we deal with datasets like this



## Some Examples of This Type of Data

• Microarray data

р	number of genes	several thousands
n	number of patients	at best a few hundreds

Image: Image:

## Some Examples of This Type of Data

• Microarray data

р	number of genes	several thousands
n	number of patients	at best a few hundreds

#### • Asthenosphere Data

р	number of locations	about 5000
n	number of days	3650 days

• Downweighting entire rows may be too "wasteful"

- Downweighting entire rows may be too "wasteful"
- Rows may be only partially spoiled

- Downweighting entire rows may be too "wasteful"
- Rows may be only partially spoiled
- Consider "cell contamination" as opposed to "row contamination"

- Downweighting entire rows may be too "wasteful"
- Rows may be only partially spoiled
- Consider "cell contamination" as opposed to "row contamination"
- Need for more flexible robustness methods and models

• From Alqallaf's PhD thesis and Alqallaf et al (2009)

 $\mathbf{X} = (I - B) \mathbf{Y} + B\mathbf{Z}$  $B = diag (B_1, B_2, ..., B_p)$  $P (B_i = 1) = 1 - P (B_i = 0) = \epsilon_i$ 

• From Alqallaf's PhD thesis and Alqallaf et al (2009)

$$\mathbf{X} = (I - B) \mathbf{Y} + B\mathbf{Z}$$
$$B = diag (B_1, B_2, ..., B_p)$$
$$P (B_i = 1) = 1 - P (B_i = 0) = \epsilon_i$$

• Lot's of room for research at the MSc and PhD levels on this area

• How to measure robustness in this more general context?

- How to measure robustness in this more general context?
- How to achieve robustness in this more genreral context?

- How to measure robustness in this more general context?
- How to achieve robustness in this more genreral context?
- Some initial work was done in Alqallaf (2009), for the multivariate location/scatter model

- How to measure robustness in this more general context?
- How to achieve robustness in this more genreral context?
- Some initial work was done in Alqallaf (2009), for the multivariate location/scatter model
- Lot's of room for research at the MSc and PhD levels on this area