

# Robust clustering

J.D. Gonzalez<sup>1</sup> V.J. Yohai<sup>2</sup> R.H. Zamar<sup>3</sup>

Universidad de Buenos Aires and University of British Columbia

- ▶ Cluster analysis is a useful tool for unsupervised data analysis.
- ▶ The goal is to divide the population in  $K$  homogeneous groups  $C_1, \dots, C_k$  called clusters.
- ▶ One way to define the clusters is by giving the centers of the clusters

$$\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$$

- ▶ Then

$$C_k = \{\mathbf{x}_i : \|\mathbf{x}_i - \boldsymbol{\mu}_k\| = \min_{1 \leq j \leq K} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|\}. \quad (1)$$

- ▶ In this talk we will discuss how to choose the centers.

- ▶ One of the most popular procedures to choose the cluster centers is called *K-means*.
- ▶ Introduced by Steinhaus (1956) and popularized MacQueen (1967).

Set

$$D(\mathbf{x}_i, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K) = \min_{1 \leq k \leq K} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|$$

Then

$$(\hat{\boldsymbol{\mu}}_1, \dots, \hat{\boldsymbol{\mu}}_K) = \arg \min_{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K} \frac{1}{n} \sum_{i=1}^n D^2(\mathbf{x}_i, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)$$

- ▶ Cuesta-Albertos et al. (1997) noted that K-means is not robust against outliers and proposed a robust alternative, called trimmed K-means.
- ▶ We propose another robust alternative called K-TAU.
- ▶ We use the concepts of scale estimators.

# Scale estimators

Given real numbers  $u_1, \dots, u_n$ , an scale estimator  $s(u_1, \dots, u_n)$  measures the absolute size of the sample.

An scale estimator has the following properties:

**P1.** Only depend on  $|u_1|, \dots, |u_n|$

**P2.**  $s(0, \dots, 0) = 0$

**P3.** If

$$|v_1| \geq |u_1|, \dots, |v_n| \geq |u_n|$$

then

$$s(v_1, \dots, v_n) \geq s(u_1, \dots, u_n)$$

**P4.**  $s(\lambda u_1, \dots, \lambda u_n) = |\lambda|s(u_1, \dots, u_n)$

# Examples of scale estimators

- ▶ L<sub>2</sub> scale

$$s_2(u_1, \dots, u_n) = \left( \frac{1}{n} \sum_{i=1}^n u_i^2 \right)^{1/2}$$

- ▶ L<sub>1</sub> scale

$$s_1(u_1, \dots, u_n) = \frac{1}{n} \sum_{i=1}^n |u_i|$$

# Examples of robust scale estimators

- ▶ Median scale

$$s_{\text{med}}(u_1, \dots, u_n) = \text{median}(|u_1|, \dots, |u_n|)$$

- ▶ M-scale  $s_M(u_1, \dots, u_n)$  solves

$$\frac{1}{n} \sum \rho \left( \frac{|u_i|}{s} \right) = b$$

$\rho : R \rightarrow [0, \infty)$  is even, non decreasing,  $\rho(\infty) = 1$

$b = 1/2$  for breakdown point equal to  $1/2$ .



# The Tau Scale

Yohai and Zamar (1988) introduce a new family of scale estimators,  $\tau$ -estimators to achieve high Gaussian efficiency and breakdown point  $1/2$  simultaneously.

To define a  $\tau$  scale estimator we use two  $\rho$  functions

$$\rho_1 \geq \rho_2$$

We first compute an M-scale  $s = s(u_1, \dots, u_n)$  using  $\rho_1$

The (squared)  $\tau$ -scale is then defined as

$$s_{\tau}^2(u_1, \dots, u_n) = s^2 \frac{1}{n} \sum_{i=1}^n \rho_2\left(\frac{u_i}{s}\right)$$

# K-Tau clustering procedure

The K-tau cluster centers are given by:

$$(\tilde{\mu}_1, \dots, \tilde{\mu}_K) = \arg \min_{\mu_1, \dots, \mu_K} s_\tau(D(\mathbf{x}_1, \mu_1, \dots, \mu_K), \dots, D(\mathbf{x}_n, \mu_1, \dots, \mu_K))$$

and the clusters  $C_1, C_2, \dots, C_K$  are defined

$$C_k = \{\mathbf{x}_i : \min_{1 \leq j \leq K} \|\mathbf{x}_i - \tilde{\mu}_j\| = \|\mathbf{x}_i - \tilde{\mu}_k\|\}, \quad (2)$$

as in the case of  $K$ -means

# Estimating equations

The cluster centers corresponding to the  $K$ -TAU clustering procedure can not be given explicitly, but they satisfy the following fixed point equations

$$\boldsymbol{\mu}_k = \frac{\sum_{i \in C_k} w\left(\frac{\|\mathbf{x}_i - \boldsymbol{\mu}_k\|}{s}\right) \mathbf{x}_i}{\sum_{i \in C_k} w\left(\frac{\|\mathbf{x}_i - \boldsymbol{\mu}_k\|}{s}\right)} \quad 1 \leq k \leq K, \quad (3)$$

with

$$w(t) = A_n \frac{\psi_1(t)}{t} + B_n \frac{\psi_2(t)}{t}. \quad (4)$$

# Estimating equations

$$A_n = \sum_{i=1}^n \left[ 2\rho_2 \left( \frac{D(\mathbf{x}_i, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)}{s} \right) - \psi_2 \left( \frac{D(\mathbf{x}_i, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)}{s} \right) \right]$$

$$B_n = \sum_{i=1}^n \psi_1 \left( \frac{D(\mathbf{x}_i, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)}{s} \right)$$

and the M-scale  $s$  satisfies

$$\frac{1}{n} \sum_{i=1}^n \rho_1 \left( \frac{D(\mathbf{x}_i, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)}{s} \right) = b$$

# Estimating equations

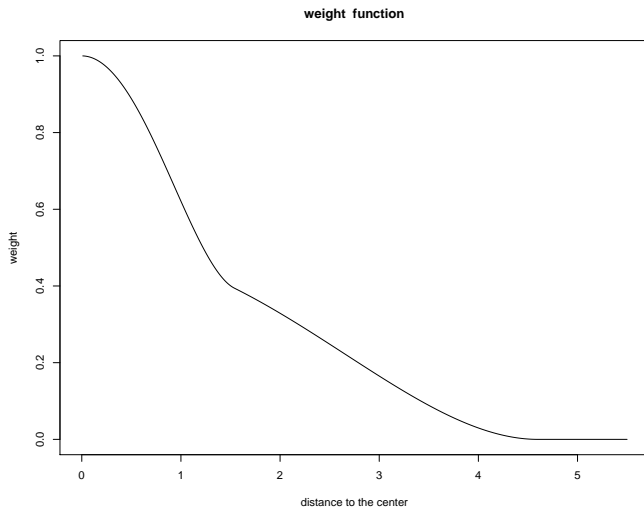


Figure: Weights based on  $\rho_1$  and  $\rho_2$  bisquare functions

The following steps are repeated  $H$  times.

**Initialization.**  $K$  initials centers  $\mu_1^{(0)}, \dots, \mu_K^{(0)}$  are obtained to start the algorithm.

**Update Centers and Groups.** Let  $\mu_1^{(l)}, \dots, \mu_K^{(l)}, C_1^{(l)}, \dots, C_K^{(l)}$ , and  $s^{(l)}$  be the current values.

The new centers are given by

$$\mu_k^{(l+1)} = \frac{\sum_{i \in C_k} w \left( \frac{\|\mathbf{x}_i - \mu_k^{(l)}\|}{s^{(l)}} \right) \mathbf{x}_i}{\sum_{i \in C_k} w \left( \frac{\|\mathbf{x}_i - \mu_k^{(l)}\|}{s^{(l)}} \right)}, 1 \leq k \leq K,$$

and  $s^{(l+1)}$  satisfies:

$$\frac{1}{n} \sum_{i=1}^n \rho_1 \left( \frac{D(\mathbf{x}_i, \mu_1^{(l)}, \dots, \mu_K^{(l)})}{s} \right) = b$$

**Stopping Rule.** Given a tolerance value  $\varepsilon$ , we stop when

$$\frac{\|\boldsymbol{\mu}_k^{(l+1)} - \boldsymbol{\mu}_k^{(l)}\|}{\|\boldsymbol{\mu}_k^{(l)}\|} \leq \varepsilon, \quad 1 \leq i \leq K.$$



# Accounting for Elliptical Shapes

For each of the clusters we compute robust S-estimators of multivariate location and covariance  $\tilde{\boldsymbol{\mu}}_k$  and  $\tilde{\boldsymbol{\Sigma}}_k$ .

The new clusters are defined by

$$C_k = \{\mathbf{x}_i : M(\mathbf{x}_i, \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k) = \min_{1 \leq j \leq K} M(\mathbf{x}_i, \tilde{\boldsymbol{\mu}}_j, \tilde{\boldsymbol{\Sigma}}_j)\}$$

with

$$M(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = [(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})]^{1/2}$$

# Identifying Outliers

Once we have the new clusters, we identify those observations that are outlying:

$$M^2(\mathbf{x}_i, \mu_k, \boldsymbol{\Sigma}_k) \geq \chi_{p,0.975}^2$$

Here  $\chi_{p,\alpha}^2$  is the  $\alpha$  quantile of the  $\chi^2$  distribution with  $p$  degrees of freedom

## Example

The following example is the set M5 of artificial data generated by García Escudero et al. (2008)

The data set consist of 2000 points from three bivariate normal populations  $N(\mu_i, \Sigma_i)$ , where

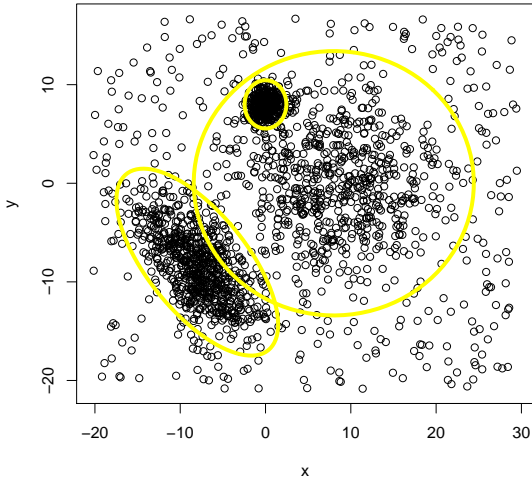
$$\mu_1 = (0, 8), \quad \mu_2 = (8, 0), \quad \mu_3 = (-8, 8)$$

$$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 45 & 0 \\ 0 & 30 \end{pmatrix}, \Sigma_3 = \begin{pmatrix} 15 & -10 \\ -10 & 15 \end{pmatrix}$$

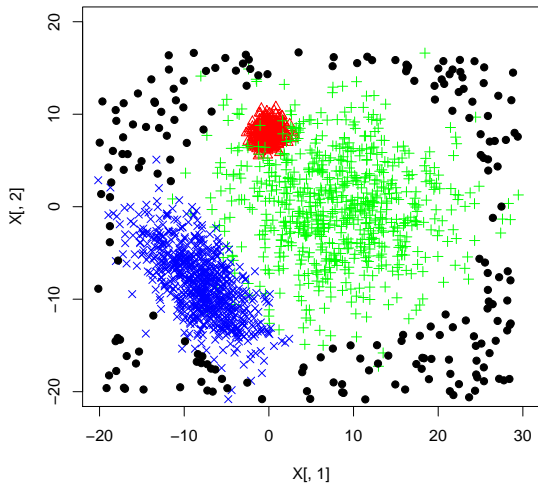
There are 1800 regular points and 200 outliers generated with a uniform distribution in a box around the three populations

From the 1800 regular points, 20% correspond to population 1, 40% to population 2 and 40% to population 3

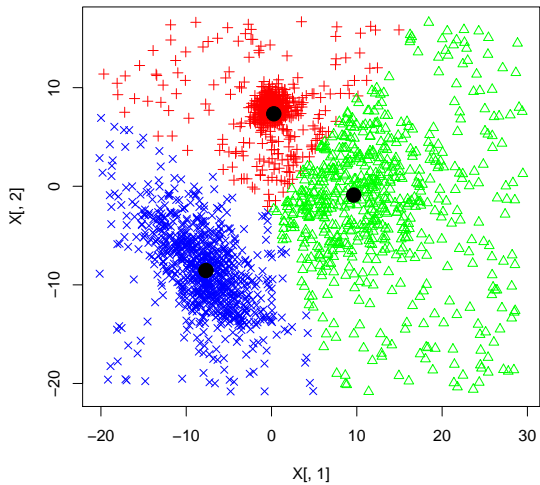
# M5Data



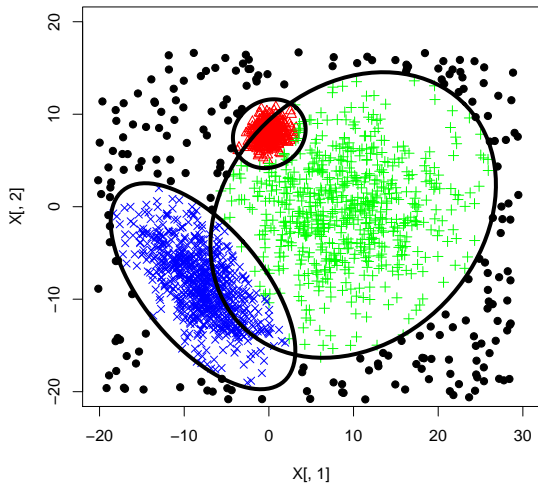
(a) M5 Data (true groups)



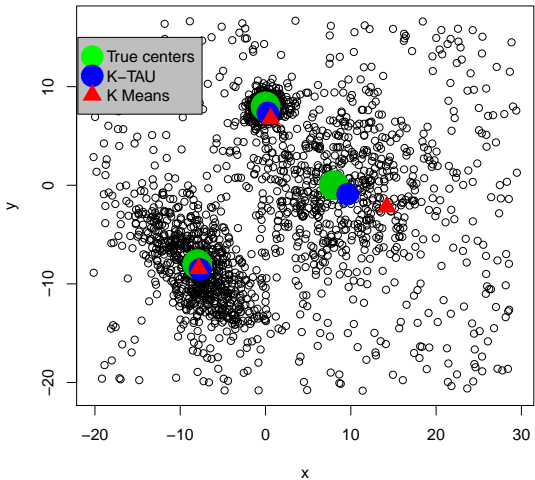
# K TAU



(b) Improved K Tau

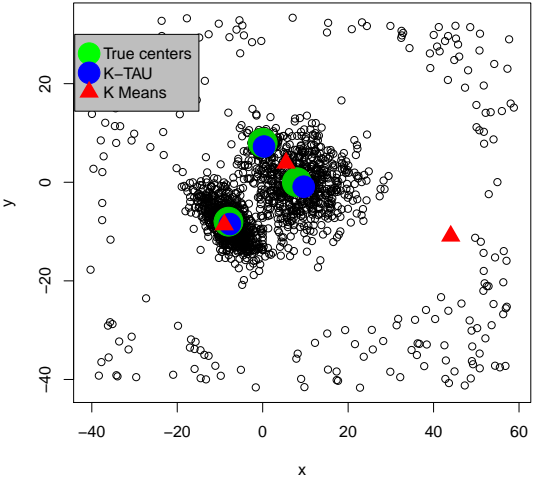


# M5 data

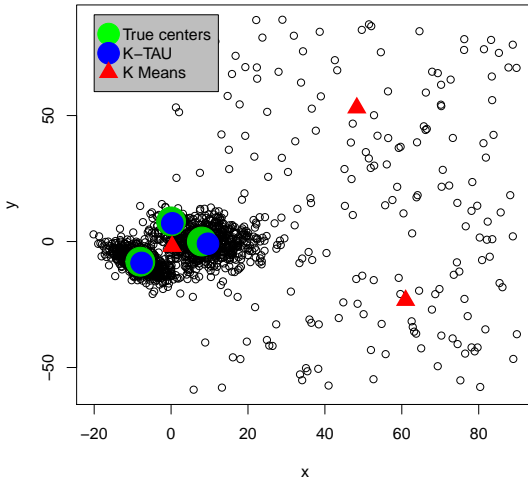




### M5 data Outliers multiplied by 2



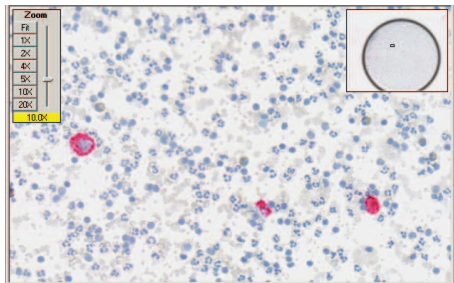
### M5 data outliers in $[0,90] \times [-60,90]$



Potential fields of application for robust cluster analysis

- ▶ Computer vision
- ▶ Anomaly detection
- ▶ Inspection of industrial items
- ▶ Search for tumors in microscopic images
- ▶ Search for a lost vessel in satellite images

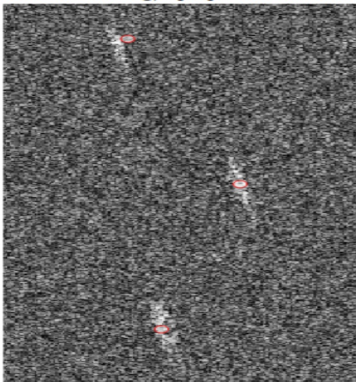
- ▶ Micrometastasis of tumor cells in circulating blood
- ▶ An algorithm tuned to find the “rare events” can process an entire slide quickly and reliably without fatigue



# Satellite Imagery Used by Frontex to Detect and Rescue Migrant Boats

While the use by Frontex of satellite imagery is not new, [Frontex released a copy of a satellite image](#) used last week to detect and rescue 370 people on board three inflatable boats off the Libyan coast. (It is unclear whether the image made available by Frontex shows the actual spatial resolution available to Frontex.)

According to Frontex, the imagery is part of "[Frontex's Eurosur Fusion Services](#) ... made possible by the cooperation between experts at Frontex and the [European Maritime Safety Agency](#) (EMSA), Italian authorities and EUNAVFORMED. ... The Eurosur [fusion] services already include automated large vessel tracking and detection capabilities, software functionalities allowing complex calculations for predicting positions and detecting suspicious activities of vessels, as well as precise weather and oceanographic forecasts. Fusion Services use optical and radar



# TELEFE Picture

TELEFE's "No-Signal" TV Image

$40 \times 50 = 2000$  pixels



## RGB-Color Coding

Each pixel corresponds to a 3-d vector:

$$( R, G, B ), \quad 0 \leq R, G, B \leq 1.$$

The vector  $( R, G, B )$  gives the intensity of red, green and blue for the pixel.

For example:

$(R, G, B) = (1, 0, 0) = \text{Red}$

$(R, G, B) = (0, 1, 0) = \text{Green}$

$(R, G, B) = (0, 0, 1) = \text{Blue}$

$(R, G, B) = (1, 1, 1) = \text{White}$

$(R, G, B) = (0, 0, 0) = \text{Black}$

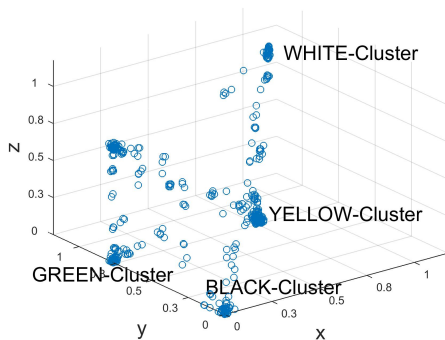
$(R, G, B) = (.5, .6, 0) = \text{Yellow}$



# TELEFE Picture Again (left side only)



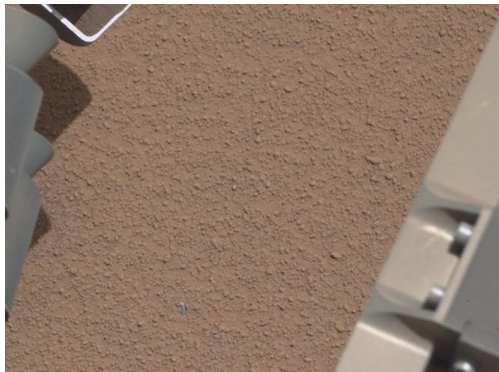
Toy Example



# Mars Rover Curiosity

Part of a high resolution NASA's picture

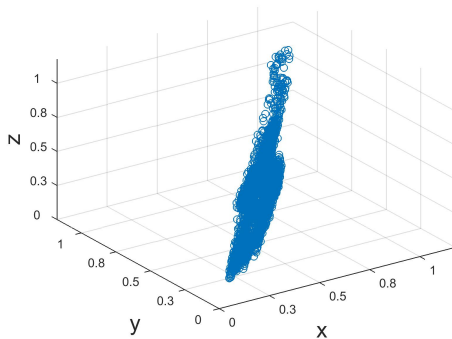
**495 × 664 = 328,680 pixels**



# Mars Rover Curiosity

RGB - Representation

**True Example**



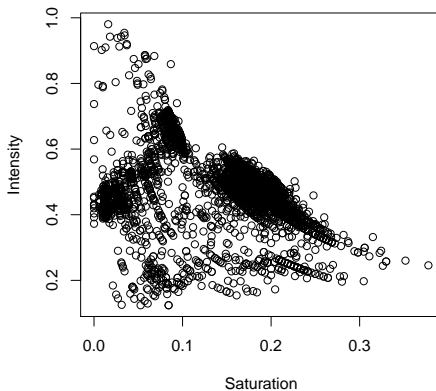
## Color picture 2-d representation

$$I = \frac{R + G + B}{3} \quad \text{Intensity}$$

$$S = 1 - \frac{3 \times \min \{R, G, B\}}{R + G + B} \quad \text{Saturation}$$

# Mars Rover Curiosity

## SI - Representation



- ▶ Image is divided into  $n = 3234$  cells with  $10 \times 10$  pixels
- ▶ Each pixel has two values  $(I, S)$
- ▶ Each observation represents a 200-d vector

$$(I_1, \dots, I_{100}, S_1, \dots, S_{100})$$

- ▶ The picture has three components (clusters):
  - ▶ *shinning metal (SHM)*
  - ▶ *opaque metal (OPM)*
  - ▶ *sand (SND)*

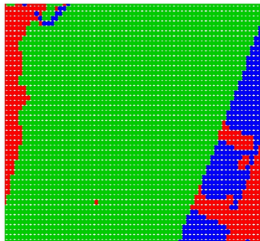
# Picture Segmentation

Robust and non-robust clustering

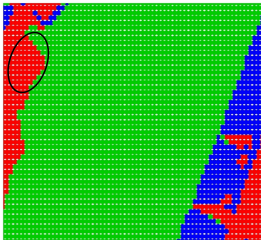
(a) Original Image



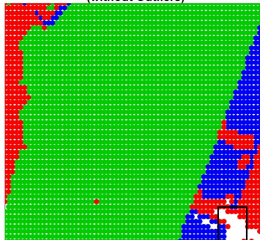
(b) TAU



(c) K Means



(d) K Means  
(without Outliers)

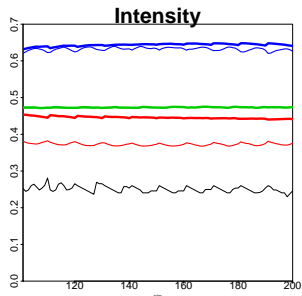
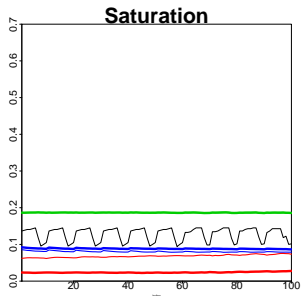


- ▶ 15 % of the OPM-cluster cells are very opaque (right lower corner in the picture)
- ▶ These cells have unusually low I-levels and high S-levels.
- ▶ These outliers upset the K-Means OPM-cluster center.



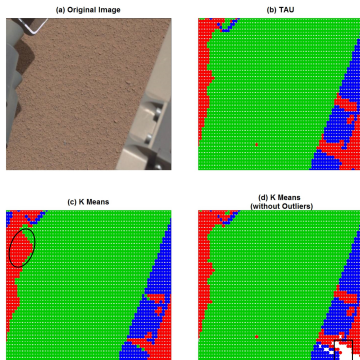


# Cluster Centers



# Effect of Outliers

- ▶ The shaded sand region is assigned to the OPM-cluster by K-means.
- ▶ Recomputing the K-means clusters after removing these outliers validates this reasoning.
- ▶ Now the K-means results are consistent with those of the robust clustering procedures.



# Automatic detection of missing objects

- ▶ The robust clustering sorted the cells as follows:
  - ▶  $n_1 = 2500$  SND-cluster cells
  - ▶  $n_2 = 405$  OPM-cluster cells
  - ▶  $n_3 = 329$  SHM-cluster cells
- ▶ The screw is made of a material – opaque metal – that makes up 13% of the image

## “Geographic” Step

- ▶ Restrict attention to the  $n_2 \times 2$  *Geographic Data Matrix* with the position (latitude and longitude) of the OPM–cluster cells
- ▶ Perform a second robust cluster analysis on these Geographic data
- ▶ Isolated outliers in this second robust clustering are likely to locate the missing piece.

# The geographic data

