Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Extreme-Value Analysis of Corrosion Data

Marcin Glegola

July 16, 2007

Supervisors:
Prof. Dr. Ir. Jan M. van Noortwijk
MSc Ir. Sebastian Kuniewski
Dr. Marco Giannitrapani

$\overset{\textit{f}}{\textbf{T}}\textbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Outline

Motivation

Objective of the Thesis

Methods used

Data declustering

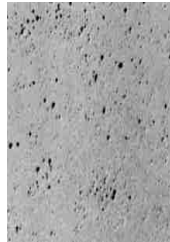Examples of application

Framework for modelling extremes of corrosion

Conclusions and recommendations

$\mathbf{\tilde{T}U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
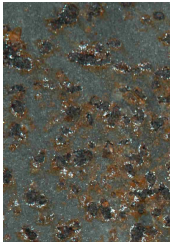Framework for modelling extremes of corrosion
Conclusions and recommendations

## Motivation

- ▶ in the oil industry, hundreds of kilometres of pipes and other equipment can be affected by corrosion

- ▶ it is extreme defect depth/wall loss that influences the system reliability ⇒ extreme-value methods are sensible for application

- ▶ only part of the system can be subjected to inspection ⇒ results extrapolation is needed

$\overset{\text{\textit{f}}}{\mathbf{T}}\mathbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Motivation

▶ defects/wall losses caused by corrosion are likely to be locally dependent ⇒ independence assumption questionable

Motivation
**Objective of the Thesis**
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
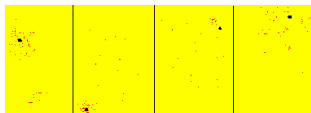Conclusions and recommendations

## Objective of the Thesis

- ▶ present statistical methods to model extremes of corrosion data, taking into account local defect dependence

- ▶ spatial extrapolation of the results

- ▶ examples of application

- ▶ framework/guideline for modelling extreme-values of corrosion

$\overset{\textit{f}}{T}U$Delft

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Methods used

The Generalised Extreme-Value (GEV) distribution (block maxima data)
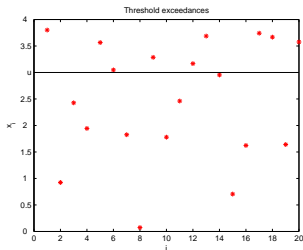


$z_1$     $z_2$     $z_3$   ...   $z_k$

$$G(z) = \begin{cases} \exp\left\{-\left[1 + \xi\left(\dfrac{z-\mu}{\sigma}\right)\right]_+^{-\frac{1}{\xi}}\right\}, & \xi \neq 0 \\ \exp\left\{-\exp\left[-\left(\dfrac{z-\mu}{\sigma}\right)\right]\right\}, & \xi = 0, \end{cases}$$

**TU**Delft

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Methods used

The Generalised-Pareto (GP) distribution (excess over threshold data)

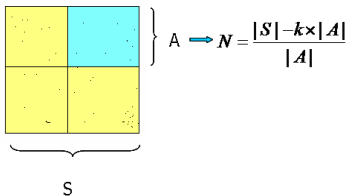$$Y_i = X_i - u, \text{ for } X > u, \ i = 1, \ldots, n_u$$



$$H(y) = \begin{cases} 1 - \left[1 + \dfrac{\xi y}{\bar{\sigma}}\right]_+^{-1/\xi}, & \xi \neq 0 \\ 1 - \exp\left(-\dfrac{y}{\bar{\sigma}}\right), & \xi = 0 \end{cases}$$

$\vec{T}$UDelft

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Methods used (extrapolation-GEV)

▶ return-level method

$$G(z_p) = 1 - p \Leftrightarrow Pr\{M > z_p\} = p = \frac{1}{N}$$



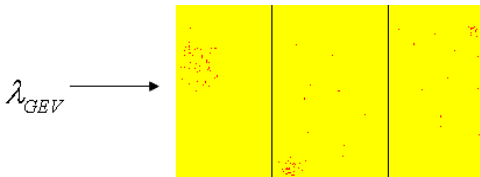$$A \longrightarrow N = \frac{|S| - k \times |A|}{|A|}$$

S

▶ implied distribution of the maximum corresponding to the not inspected area

$$Pr\{X_N \leq z\} = G_N(z) = G(z)^N$$

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Methods used (extrapolation-GP)

based on Poisson frequency of threshold exceedances (Poisson-GP model)
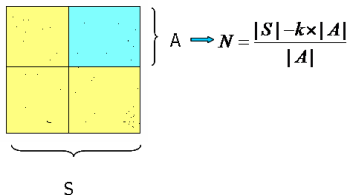
$\lambda_{GEV} \longrightarrow$ 

▶ return-level method

$$H(y_p) = 1 - p \Leftrightarrow Pr\{Y > y_p\} = p = \frac{1}{N_E}$$

$N_E = \lambda_{GEV} \times (|S| - k \times |A|)$ - expected number of exceedances on the not inspected area

$\tilde{T}$UDelft

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Methods used (extrapolation-GP)

▶ implied distribution of the maximum corresponding to the not
inspected area

$$Pr\{X_N \leq x\} = \exp\left\{-\lambda_{GEV}\left(1 + \xi\frac{x - u}{\bar{\sigma}}\right)_+^{-1/\xi}\right\}^N$$



$$A \Longrightarrow N = \frac{|S| - k \times |A|}{|A|}$$

S

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations
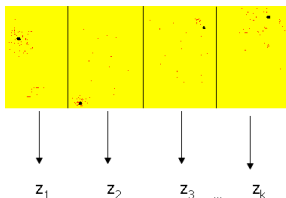
## Methods used-summary

- ▶ two methods for statistical inference about extreme-values of corrosion
  - the GEV distribution (block maxima)
  - the GP distribution (excess over threshold data)

- ▶ two methods for spatial results extrapolation
  - return-level
  - distribution of the maximum corresponding to the not inspected area

- ▶ the GEV and GP distributions are closely related and theoretically, should give the same results

$\widetilde{T}U$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Modelling extremes of dependent data

- ▶ for the stationary data characterised by the limited extend of long-range dependence at extreme levels, the extreme-value methods can be still applied

- ▶ in corrosion application it is reasonable to assume that pit depths are locally dependent $\Rightarrow$ extreme-value methods are applicable

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
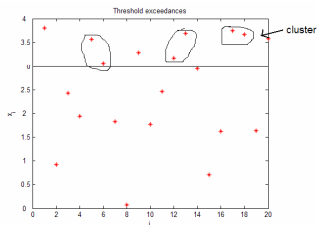Conclusions and recommendations

# Modelling extremes of dependent data with the **GEV** distribution

- ▶ assuming local data dependence, block maxima of stationary data (for sufficiently large block sizes) can be considered as approximately independent

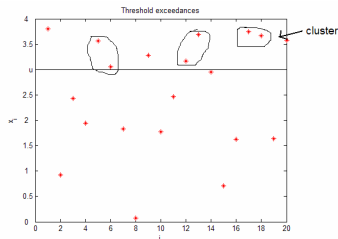- ▶ the GEV distribution is used in its standard form



$z_1$  $z_2$  $z_3$  ...  $z_k$

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Modelling extremes of dependent data with the **GP** distribution

- ▶ neighbouring exceedances may be dependent, therefore the change of practise is needed

- ▶ one of the most widely adopted method is **data declustering** - filtering out dependent observations such that remaining exceedances can be considered as approximately independent
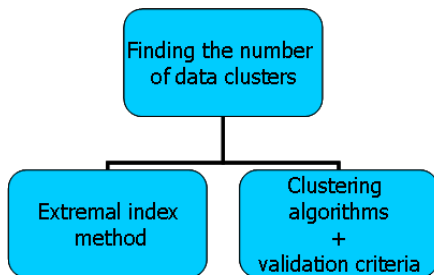


$\vec{T}$UDelft

Motivation
Objective of the Thesis
**Methods used**
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Modelling extremes of dependent data with the **GP** distribution-approach

▶ define clusters of exceedances

▶ identify maximum excess within each cluster

▶ assuming that cluster maxima are independent fit the GP distribution



$\tilde{\mathbf{T}}$UDelft

# Data declustering

Estimation of the number of data clusters

Motivation
Objective of the Thesis
Methods used
**Data declustering**
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

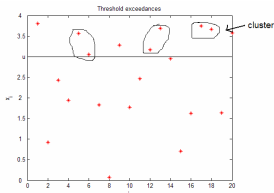## Data declustering

Extremal index method:

- ▶ the extent of short-range dependence of extreme events is captured by the parameter $\theta$, called **extremal index**

$$\theta = \frac{1}{\text{limiting mean cluster size}}$$

- ▶ extremal index measures the degree of clustering of the process at extreme levels

$\overset{\text{\textit{f}}}{\textbf{T}}\textbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Data declustering

Extremal index method



$$\theta = \frac{1}{\text{limiting mean cluster size}} \Rightarrow N_c = \theta \times N_e$$

where $N_c$ - number of clusters, $N_e$ - number of exceedances above threshold u
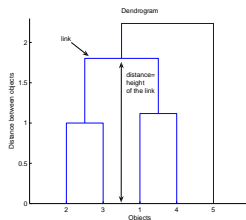
$\theta$ is estimated by the intervals estimator

$\overset{\text{\textit{f}}}{\textbf{T}}\textbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
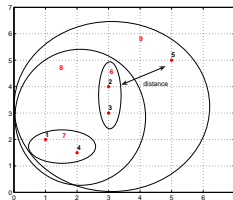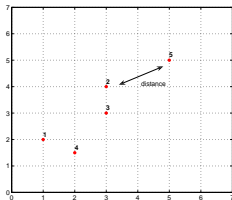Conclusions and recommendations

## Data declustering

Finding number of clusters using clustering algorithm

- ▶ define a validity criteria for the number $N_c$ of found clusters

- ▶ run the clustering algorithm for a range of $N_c$

- ▶ as proper number of clusters choose the one for which the validity criteria are optimised

$\tilde{T}$UDelft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Agglomerative hierarchical clustering algorithm

- ▶ starts with single points as clusters

- ▶ at each step the two closest clusters are merged

- ▶ stops when only one cluster remains

$\mathbf{\tilde{T}U}$Delft

Motivation
Objective of the Thesis
Methods used
**Data declustering**
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Agglomerative hierarchical clustering algorithm

Motivation
Objective of the Thesis
Methods used
**Data declustering**
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Data declustering

Validation criteria, that we used, aim at identifying clusters
that are compact and well isolated:

► silhouette plot - maximum value indicates optimum

► Davies-Bouldin index - minimum indicates optimum

$\overset{\prime}{\textbf{T}}\textbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Data declustering-summary

- ▶ two methods to estimate the number of data clusters
  - the extremal index method
  - clustering algorithm $+$ validation criteria (Davies-Bouldin index, silhouette plot)

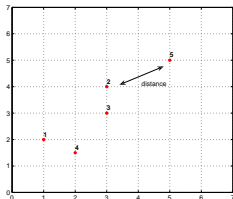- ▶ prior to data declustering perform clustering tendency test

$\overset{\textit{f}}{\textbf{T}}\textbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

## Simulated corroded surface

▶ application of the gamma-process model
▶ dependence in terms of the product moment correlation

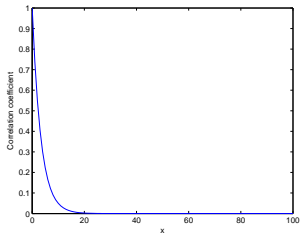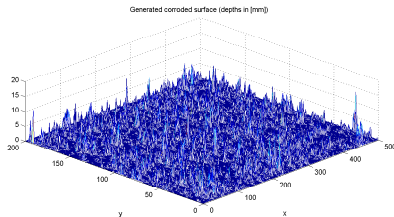$$\rho(X_k, X_l) = \exp \left\{ -d \left( \sum_{i=1}^{2} |dist_i|^p \right)^{q/p} \right\}$$

$X_k = (x_k, y_k), \ X_l = (x_l, y_l),$
$dist_1 = |x_k - x_l|, \ dist_2 = |y_k - y_l|,$
$d = 0.3, \ p = 2, \ q = 1$



$\vec{T}$UDelft

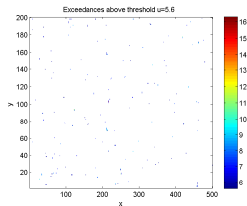Motivation
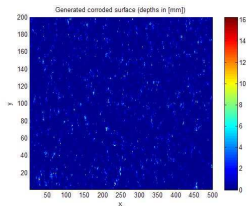Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Simulated corroded surface

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Simulated corroded surface

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Simulated corroded surface - clustering algorithm

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Simulated corroded surface - results

| $\hat{\theta}$ | $n_c$ | $n_{c_{alg}}$ |
|---|---|---|
| 0.544 | 107 | 121 |

Table: The estimate of extremal index and determined number of clusters

| Number of clusters | $AD^2_{up}$ $p-v.$ | KS $p-v.$ |
|---|---|---|
| 107 | 0.376 | 0.204 |
| 121 | 0.549 | 0.493 |

Table: Goodness-of-fit test results for different number of clusters

$\overset{\textit{é}}{T}U$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

| $\hat{\xi}$ | $\hat{\bar{\sigma}}$ | $AD_{up}^2$ $p-v.$ | $KS$ $p-v.$ |
|---|---|---|---|
| 0.065 | 1.341 | 0.647 | 0.488 |

Table: GP fit to excess of dependent data

| $\hat{\xi}$ | $\hat{\sigma}$ | $AD_{up}^2$ $p-v.$ | $KS$ $p-v.$ |
|---|---|---|---|
| -0.0137 | 1.6839 | 0.555 | 0.493 |

Table: GP fit to excess of declustered data

| $\hat{\xi}$ | $\hat{\sigma}$ | $\hat{\mu}$ | $AD_{up}^2$ $p-v.$ | $KS$ $p-v.$ |
|---|---|---|---|---|
| -0.007 | 1.627 | 5.637 | 0.621 | 0.899 |

Table: GEV fit to block maxima data

$\tilde{T}U$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Real data example





half pipe

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Real data example

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

# Real data example

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

| $\hat{\xi}$ | $\hat{\sigma}$ | $\hat{\mu}$ | $AD_{up}^2\ p-v.$ | $KS\ p-v.$ |
|---|---|---|---|---|
| -0.082 | 0.182 | 0.757 | 0.713 | 0.986 |

Table: GEV fit to block maxima data

| $\hat{\xi}$ | $\hat{\sigma}$ | $AD_{up}^2\ p-v.$ | $KS\ p-v.$ |
|---|---|---|---|
| -0.008 | 0.133 | 0.616 | 0.074 |

Table: GP fit to excess of dependent data

| $\hat{\xi}$ | $\hat{\sigma}$ | $AD_{up}^2\ p-v.$ | $KS\ p-v.$ |
|---|---|---|---|
| -0.069 | 0.172 | 0.717 | 0.654 |

Table: GP fit to excess of declustered data

$\vec{\mathbf{T}}\mathbf{U}$Delft

Motivation
Objective of the Thesis
Methods used
Data declustering
**Examples of application**
Framework for modelling extremes of corrosion
Conclusions and recommendations

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
Framework for modelling extremes of corrosion
Conclusions and recommendations

Framework for modelling extremes of corrosion with the GEV distribution

Motivation
Objective of the Thesis
Methods used
Data declustering
Examples of application
**Framework for modelling extremes of corrosion**
Conclusions and recommendations

Framework for modelling extremes of corrosion with the GP distribution

Motivation
Objective of the Thesis
Methods used
Data declustering
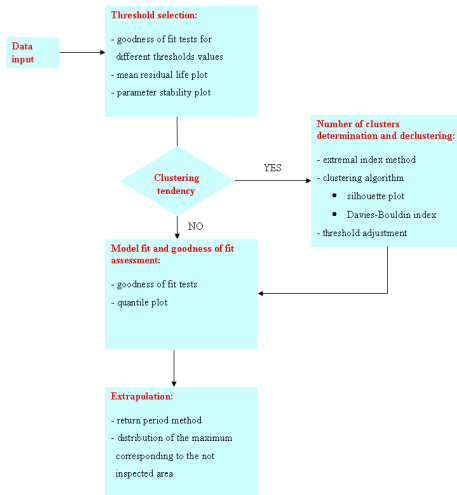Examples of application
Framework for modelling extremes of corrosion
**Conclusions and recommendations**

## Conclusions and recommendations

- ▶ the two applied distributions are closely related and lead to the consistent inference about extreme-values of corrosion

- ▶ data declustering improves the results given by the GP distribution

- ▶ the performance of other clustering algorithms could be checked

- ▶ in order to take into account corrosion nonstationarity due to space-varying environmental conditions, covariate-dependent extreme-value models with trends could be considered

$\tilde{T}$UDelft

# THANK YOU FOR ATTENTION

## Questions???

T U Delft