

# Anomaly Detection

Day 1: What is an anomaly? How do we spot them?

Day 2: Machine learning, clustering and classification

Day 3: Periodicity, and an entity centric view



Peter Tillotson

Acumed Consulting

June 5, 2016

Module Aims

Overview

What is Anomaly Detection?

What is normal?

What methods are used?

Linear regression

Classification

Clustering

Python Pandas

Data structures

Sequences

DataFrames

IO: loading and writing data

Plotting graphs

Linear Regression

Exercise 1: Anomaly detection by  
linear regression in Pandas

# Day I

What is an anomaly? How do we spot  
them?

## Module Aims

### Overview

What is Anomaly Detection?

What is normal?

What methods are used?

Linear regression

Classification

Clustering

### Python Pandas

Data structures

Sequences

DataFrames

IO: loading and writing data

Plotting graphs

### Linear Regression

Exercise 1: Anomaly detection by  
linear regression in Pandas

## 1 Module Aims

## 2 Overview

## 3 Python Pandas

## 4 Linear Regression

## 5 Exercise 1: Anomaly detection by linear regression in Pandas

## Module Aims

## Overview

What is Anomaly Detection?

What is normal?

What methods are used?

Linear regression

Classification

Clustering

## Python Pandas

Data structures

Sequences

DataFrames

IO: loading and writing data

Plotting graphs

## Linear Regression

Exercise 1: Anomaly detection by  
linear regression in Pandas

By the end of today, we will have:

- a good understanding of anomaly detection
- discussed in detail a number of techniques commonly used to detect anomalies
  - linear regression
  - classification
  - clustering
- introduced Python Panda's and
- implemented a linear regression model in with Panda's

## Module Aims

## Overview

What is Anomaly Detection?

What is normal?

What methods are used?

Linear regression

Classification

Clustering

## Python Pandas

Data structures

Sequences

DataFrames

IO: loading and writing data

Plotting graphs

## Linear Regression

Exercise 1: Anomaly detection by  
linear regression in Pandas

<b>Programming:</b> basic understanding of programming concepts including: <ul style="list-style-type: none"><li>• flow control</li><li>• functions</li><li>• data structures</li><li>• types (int, float)</li></ul>	Good
<b>Python:</b> be familiar with the Python	Some
<b>Math:</b> basic statistics (mean, variance)	Some

## Module Aims

## Overview

## What is Anomaly Detection?

What is normal?

What methods are used?

Linear regression

Classification

Clustering

## Python Pandas

Data structures

Sequences

DataFrames

IO: loading and writing data

Plotting graphs

## Linear Regression

Exercise 1: Anomaly detection by  
linear regression in Pandas

Anomaly from the Oxford Dictionary:

*Something that deviates from what is standard, normal, or expected*

But that includes 5% of all data so what we're really interested in are events that are particularly odd or inherently risky. Ideally something that is actionable and that probably needs human intervention.

We want to avoid false positives, though whilst potentially rare these create noise and reduce an analysts trust of an automated solution.

We also want to avoid false negatives, real risks that we failed to spot.

## Module Aims

## Overview

What is Anomaly Detection?

**What is normal?**

What methods are used?

Linear regression

Classification

Clustering

## Python Pandas

Data structures

Sequences

DataFrames

IO: loading and writing data

Plotting graphs

## Linear Regression

Exercise 1: Anomaly detection by  
linear regression in Pandas

```
1 import json
   def my_function():
3     m_str = json.dumps({'name': 'peter', 'msg': 'dont it look
       nice'})
   print
5     my_function()
```

**[1]**

More often than not, time series data are 'non-stationary'; that is, the values of the time series do not fluctuate around a constant mean or



with a constant variance.

## Day II

# Machine learning, clustering and classification

## Day III

# Periodicity, and an entity centric view

Anything that repeats with a fixed interval can be said to be periodic. If you can be guaranteed delivery of your signal, and have a trusted source for timestamps then it is possible to run Discrete Fourier Transforms and get a frequency domain representation of your signal.

Often this is not possible and the following approaches can be used to handle trends and seasonality in data.

- Curve fitting - STL / Decomposition
- Twitter's Seasonal Hybrid ESD [2]
  - Time series decomposition
  - Generalised ESD
- Numenta's Hierarchical Temporal Memory [3]

# Generalized ESD (extreme Studentized deviate) I

The Generalized ESD[4] test is defined for the two hypothesis:

$H_0$       There are no outliers in the data set

$H_r$       There are up to  $r$  outliers in the data set

Compute:

$$R_i = \frac{\max_j |x_j - \mu|}{\sigma} \quad (1)$$

with  $\mu$  and  $\sigma$  denoting the mean and standard deviation, respectively.

Remove the observation that maximizes  $|x_j - \mu|$  and then recompute the above statistic with  $n - 1$  observations. Repeat this process until  $r$  observations have been removed. This results in the  $r$  test statistics  $R_1, R_2, \dots, R_r$ .

## Generalized ESD (extreme Studentized deviate) II

Corresponding to the  $r$  test statistics, compute the following  $r$  critical values:

$$\lambda_i = \frac{(n-1)t_{p,n-i-1}}{\sqrt{(n-i-1+t_{p,n-i-1}^2)(n-i+1)}} \quad (2)$$

for  $i = 1, 2, \dots, r$

where  $t_{p,v}$  is the  $100p$  percentage point from the  $t$  distribution with  $v$  degrees of freedom and

$$p = 1 - \frac{\alpha}{2(n-i+1)}$$

The number of outliers is the largest  $i$  such that  $R_i > \lambda_i$ .

# REFERENCES

### References

- [1] Y. B. Nikolay Laptev Saeed Amizadeh. (Mar. 2015), A benchmark dataset for time series anomaly detection, Yahoo, [Online]. Available: <http://yahoolabs.tumblr.com/post/114590420346/a-benchmark-dataset-for-time-series-anomaly>.
- [2] A. Kejariwal. (Jan. 2015), Introducing practical and robust anomaly detection in a time series, [Online]. Available: <https://blog.twitter.com/2015/introducing-practical-and-robust-anomaly-detection-in-a-time-series>.
- [3] J. Hawkins. (Sep. 2011), Hierarchical temporal memory, [Online]. Available: <http://numenta.com/assets/pdf/whitepapers/hierarchical-temporal-memory-cortical-learning-algorithm-0.2.1-en.pdf>.
- [4] B. Rosner, "Percentage points for a generalized esd," *Technometrics*, vol. 25, no. 2, pp. 165–172, May 1983.