

# Anomaly Detection

Robert Winslow

GalvanizeU

July 13th, 2017

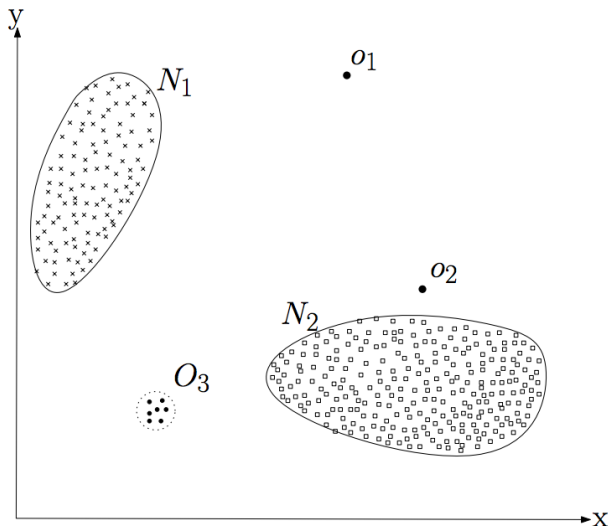
# Anomaly detection

What is an anomaly?

Anomalies are patterns in data that are not “normal”.

# Anomaly detection

What is an anomaly?



# Anomaly detection

The task of anomaly detection is to determine if an instance is *typical* or *anomalous*.

*normal* vs *abnormal*

*expected* vs *unexpected*

An anomaly is *interesting* to the analyst.

# Anomaly detection

Example

Credit card fraud

Banks lose 11 billion USD per year.

Customers lose 5 billion USD per year.

How can credit card companies deny fraudulent transactions?

... Without denying valid transactions?

# Anomaly detection

Example

Medical diagnosis

Given a brain scan. . .

Is the patient likely healthy or not?

Challenge: patient can be unhealthy in novel ways.

We cannot just build a classifier for healthy/not-healthy.

# Anomaly detection

## Example

Intrusion detection in computer security

What is typical user behavior?

Ports open to the internet on a college campus.

Is a hosted service intentional, or the result of a computer virus?

# Anomaly detection

Example

Distributed denial of service attack

Did we just have a piece of content go viral?

Or, are we victims of a botnet flooding our servers with traffic?

(Is it possible to tell the difference?)



# Anomaly detection

Example

Science experiments

LIGO: Laser Interferometer Gravitational-Wave Observatory

Identifying rare events: measurable gravity waves.

Finding needles in the cosmic haystack.

# Anomaly detection

Example

Mechanical systems monitoring.

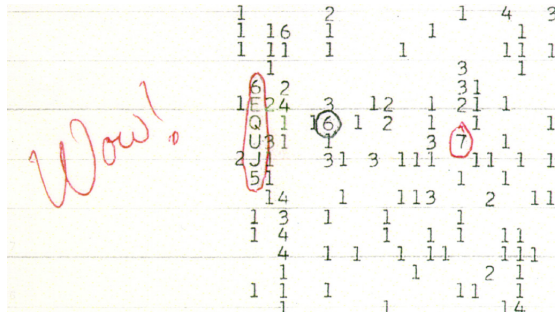
LHC cryo page

# Anomaly detection

## Example

# Detecting abnormal radio signals from space.

Wow! signal.



ET?

# Anomaly detection

Why is anomaly detection hard?

Anomaly detection would be easy if we had a training set of all anomalies.

# Anomaly detection

## Challenge

Defining the “normal” region that encompasses every possible “normal” behavior is very difficult.

# Anomaly detection

## Challenge

The border between “normal” and “abnormal” is often imprecise.

# Anomaly detection

## Challenge

If anomalies are the result of malicious actors, they can (and will) adapt to your actions.

# Anomaly detection

## Challenge

Modeling error / too little information

Example: If two customers have the same recent purchase history (both buy a BART card this morning), and one of them is a fraudulent purchase and the other is not, how could we tell the difference?



# Anomaly detection

## Challenge

The definition of “typical” or “normal” keeps changing.

Did a parameter in the underlying generative process change?

(Was today hot because of increasing global temperatures, or just the typical daily variation?)

# Anomaly detection

Challenge

Domain-dependence

Techniques are often heavily reliant on features of a particular problem domain.

Medical diagnosis: very sensitive, we look for any abnormalities.

Stock prediction: naturally high variance, so we use low sensitivity.

# Anomaly detection

## Challenge

We can almost never train a classifier for “normal” vs “abnormal”: lack of training data.

# Anomaly detection

Challenge

Data often contains noise.

Is the data point an anomaly?

Or just some Gaussian-distributed  $\epsilon$ ?

# Anomaly detection

Types of anomalies

Point anomalies

If a single data point can be judged based on its similarity to the rest of the dataset.

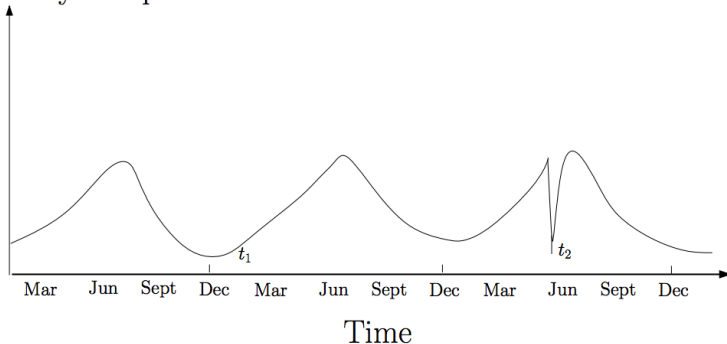
# Anomaly detection

Types of anomalies

Contextual anomalies

If a single data point is atypical within its particular context.

Monthly Temp

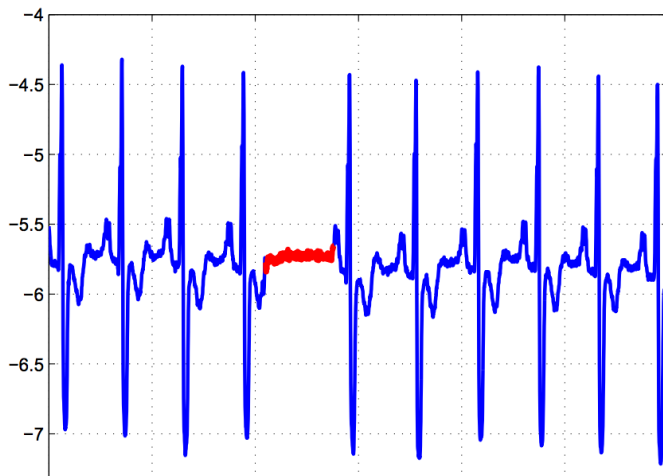


# Anomaly detection

Types of anomalies

Collective anomalies

A data point is atypical given its surrounding data points.



# Anomaly detection

Models

Supervised learning

Build a classifier: anomalous or not?

Challenges?



# Anomaly detection

Models

Semi-Supervised learning

Only have training data for typical points, not atypical ones.

Think: if data points are in low-probability regions.

Challenges?

# Anomaly detection

Models

Unsupervised learning

Think: dimensionality reduction like clustering

Challenges?

# Anomaly detection

## Outputs

Score: how anomalous is this instance?

A human or other system can then decide whether to take action.

# Anomaly detection

## Outputs

Label: is this instance typical or atypical?

Opaque, but perhaps easier to create a model for this.

Example: naive Bayes can give a probability output, but it is a bad estimate of the true probability.

But, NB can be fine if just using it for labeling.

# Anomaly detection

Generative model interpretation:

Is this data point “close” to our training data?

How likely is it that this instance can be generated from our model?

Given a model  $f$  with parameters  $\sigma \dots$

$$P(y|f; \sigma)$$

# Anomaly detection

Connections

“Happy families are all alike; every unhappy family is unhappy in its own way.” -Tolstoy

Manifold learning

Dog or cat?

Dog or cat or *other*?