

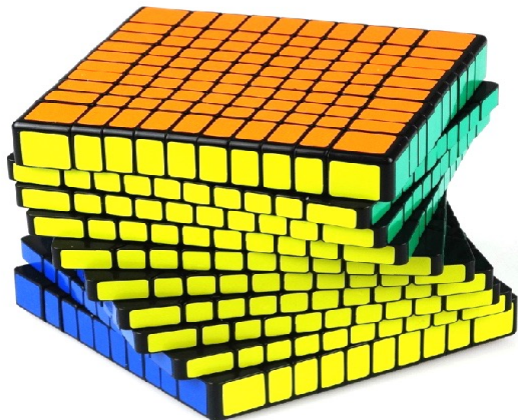
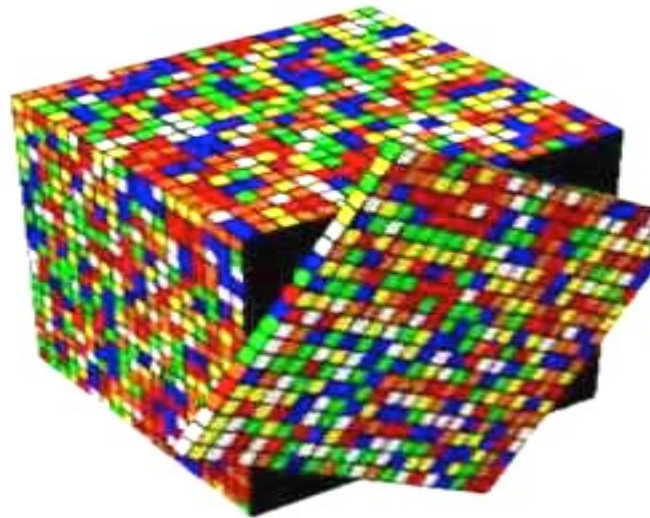
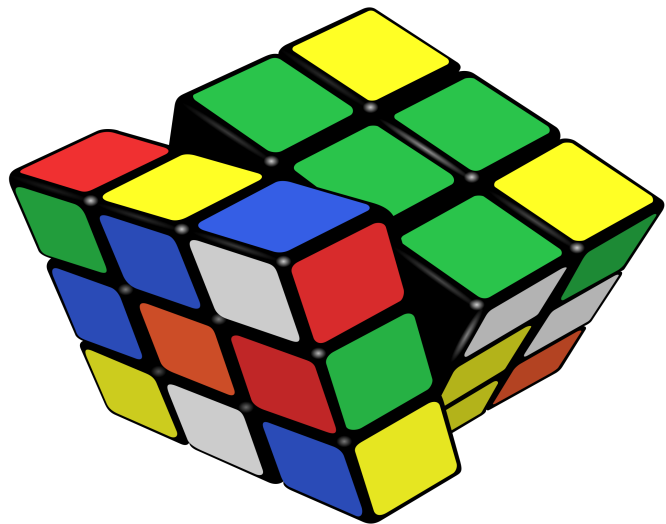
# Data Mining in Astronomy

Nadeem Oozeer, PhD

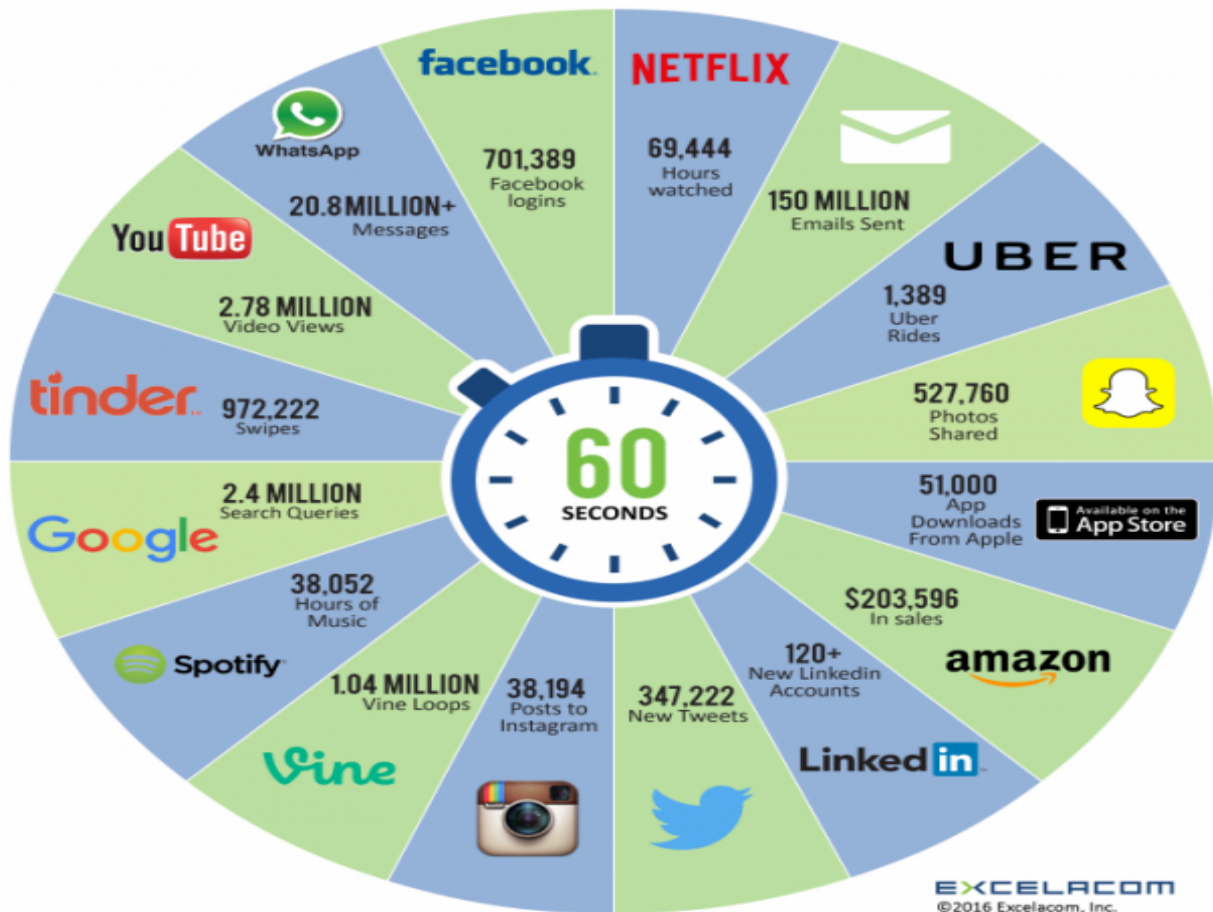
Data Scientist

SKA SA / AIMS





# 2016 What happens in an INTERNET MINUTE?



# Astronomical Data Deluge



## Square Kilometre Array



€1.5b

+ A €1.5 billion global science project



+ Astronomers and engineers from more than 70 institutes in 20 countries



3000

+ 3000 dishes, each 15m wide



+ Using enough optical fibre to wrap twice around the Earth

1,000,000 m<sup>2</sup>

+ A combined collecting area of about one square kilometre



In excess of 1 Exabyte of raw data in a single day - more than the entire daily internet traffic

Megadata



Enough raw data to fill over 15 million 64GB iPods every day



- + Automated data classification = faster with fewer errors
- + Guided search = easier access for scientists and non-scientists alike
- + Frees researchers to be more productive and creative



IBM  
Information  
Intensive  
Framework

A prototype software architecture to manage the megadata generated by SKA

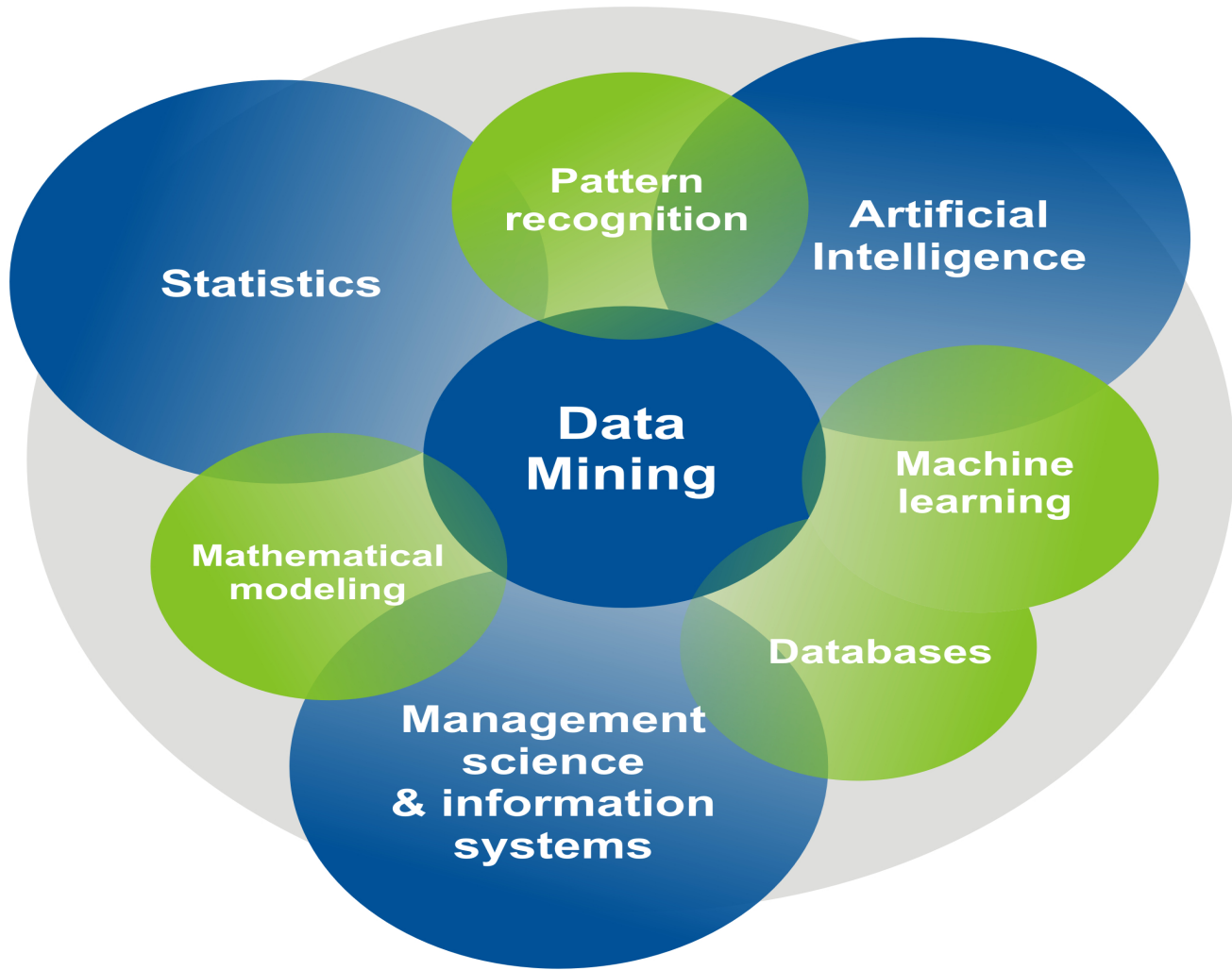


# Big Data Examples



## BIG DATA IN INDUSTRY





**Statistics**

**Pattern  
recognition**

**Artificial  
Intelligence**

**Data  
Mining**

**Machine  
learning**

**Mathematical  
modeling**

**Databases**

**Management  
science  
& information  
systems**





# Data Mining & Privacy

- Privacy
  - I want information to be used only for my benefit
- Confidentiality:
  - I want information to go only to those authorized
- Cryptography community understands confidentiality
  - Solid, vetted definitions – Proof techniques
- Not sure if anybody really understands privacy
  - *But confidentiality often sufficient*

# Solution

- Data Obfuscation
  - Nobody see the real data
- Summarization
  - Only the needed facts are exposed
- Data Separation
  - Data remains with trusted parties

**Can we afford this?**

## Illuminating the black box

Note to biologists: submissions to *Nature* should contain complete descriptions of materials and reagents used.

This journal aims to publish papers that are not only interesting and thought-provoking, but reproducible and useful. In order to do this, novel materials and reagents need to be carefully described and readily available to interested scientists.

That might seem obvious. But despite the efforts of our editors and referees, papers in the biological sciences are still being submitted — and occasionally published — that do not adequately describe the reagents used. Unless efforts are redoubled to eliminate this

established didn't want the author to reveal the sequences, as this would jeopardize its *raison d'être*. This kind of stalemate matters, because it prevents the replication of experiments and inhibits the selection of appropriate controls in subsequent work.

Some authors claim replication is possible without full sequence information or the details of novel compounds. They say that the materials in question are for sale, enabling anyone to duplicate the paper. This misses the point. Scientific progress revolves around pro-

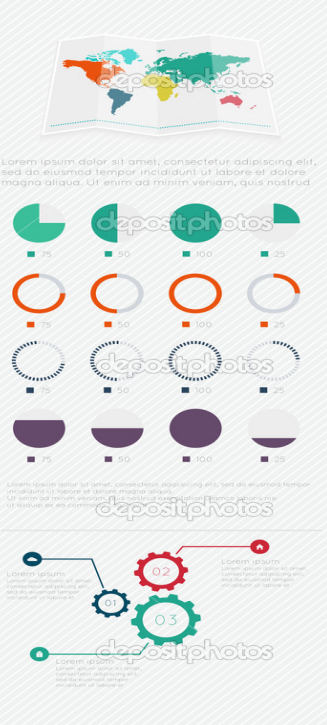
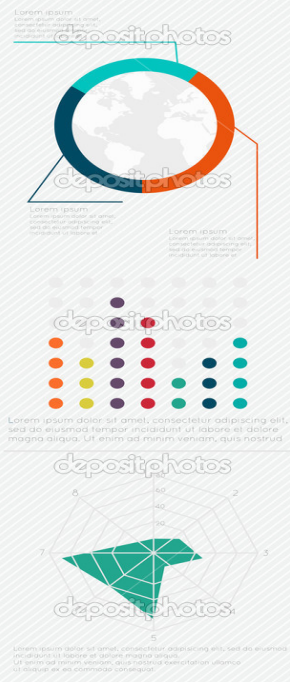
Reproducible Science  
means context, quality, trust  
means easy access to the sources

# Knowledge Discovery



# Statistics, Data Mining & Machine Learning

## INFOGRAPHICS



# Statistics

<b><i>Statistic</i></b>	<b><i>Formula</i></b>	<b><i>Used For</i></b>
Sample mean (average)	$\bar{x} = \frac{\sum x}{n}$	Measure of center; affected by outliers
Median	$n$ odd: middle value of ordered data  $n$ even: average of the two middle values	Measure of center; not affected by outliers
Sample standard deviation	$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$	Measure of variation; "average" distance from the mean
Correlation coefficient	$r = \frac{1}{n - 1} \sum \frac{(x - \bar{x})(y - \bar{y})}{s_x s_y}$	Strength and direction of linear relationship between X and Y

The Posterior

The Evidence

The Prior

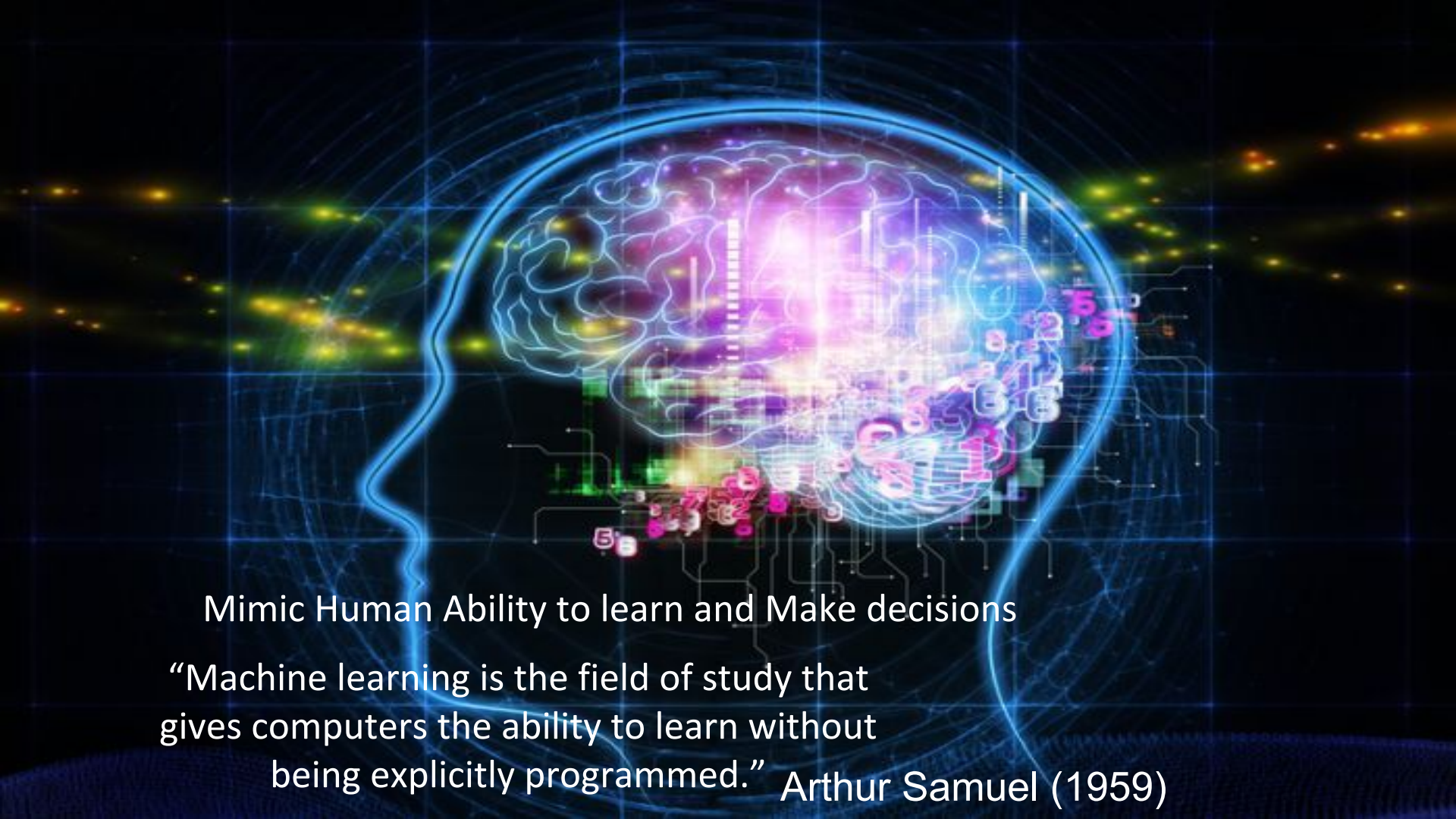
The probability of getting this evidence if this hypothesis were true

The probability of H being true, before gathering evidence

$$P(H|E) = \frac{P(H|E) P(H)}{P(E)}$$

The probability that the hypothesis (H) is true given the evidence (E)

The marginal probability of the evidence (Prob of E over all possibilities)



Mimic Human Ability to learn and Make decisions

“Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.” Arthur Samuel (1959)



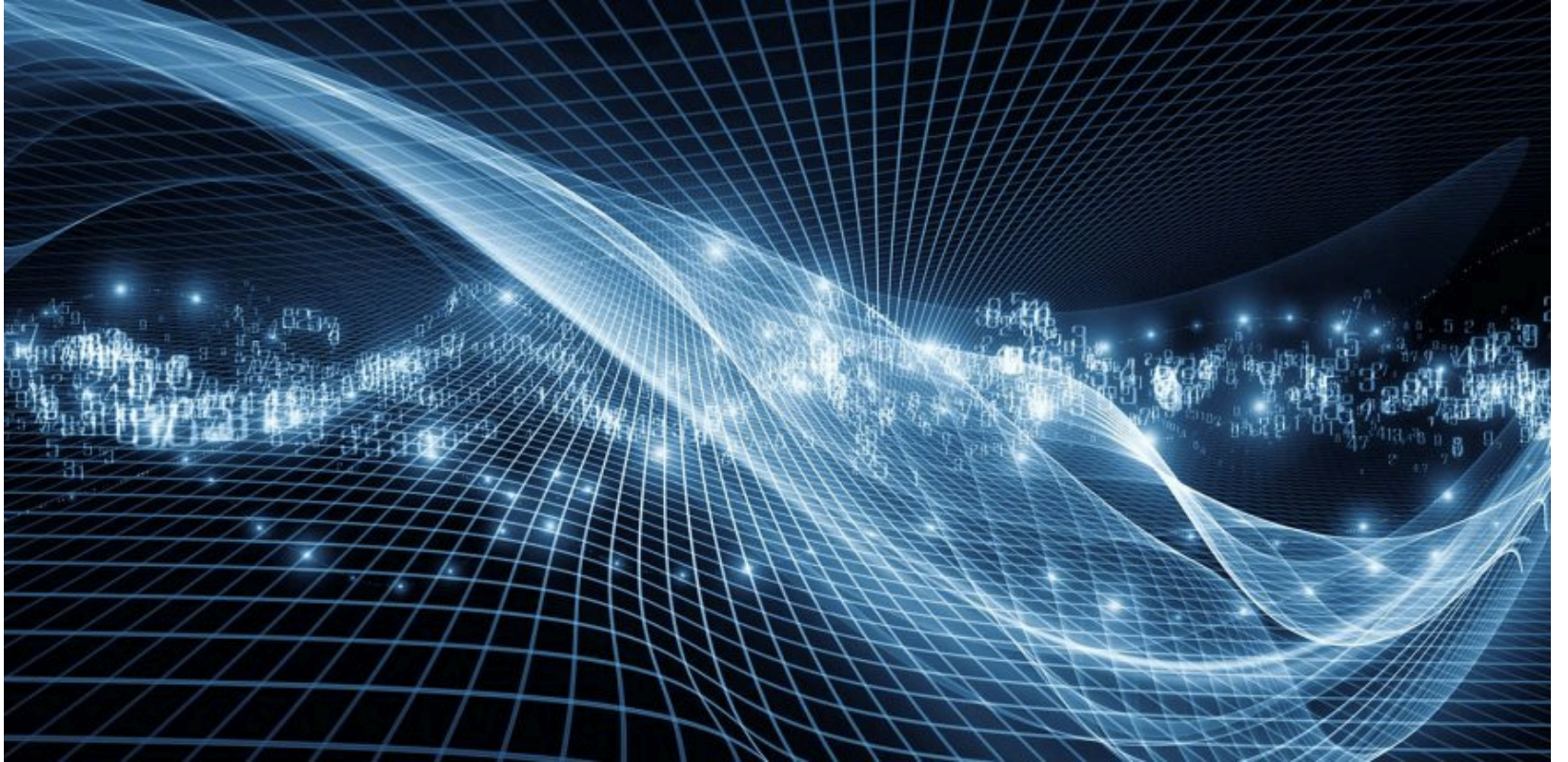
# Machine learning:

- **Supervised vs Unsupervised.**
  - *Supervised learning* - the presence of the outcome variable is available to guide the learning process.
    - there **must** be a training data set in which the solution is already known.
  - *Unsupervised learning* - the outcomes are unknown.
    - cluster the data to reveal meaningful partitions and hierarchies

# Why Data Mining?

- Classification
- Clustering
- Associations
- Visualization
- Summarization
- Serendipity ...

# When is Data enough?





# Clustering:

- Clustering is the task of gathering samples into groups of similar samples according to some predefined similarity or dissimilarity measure



sample

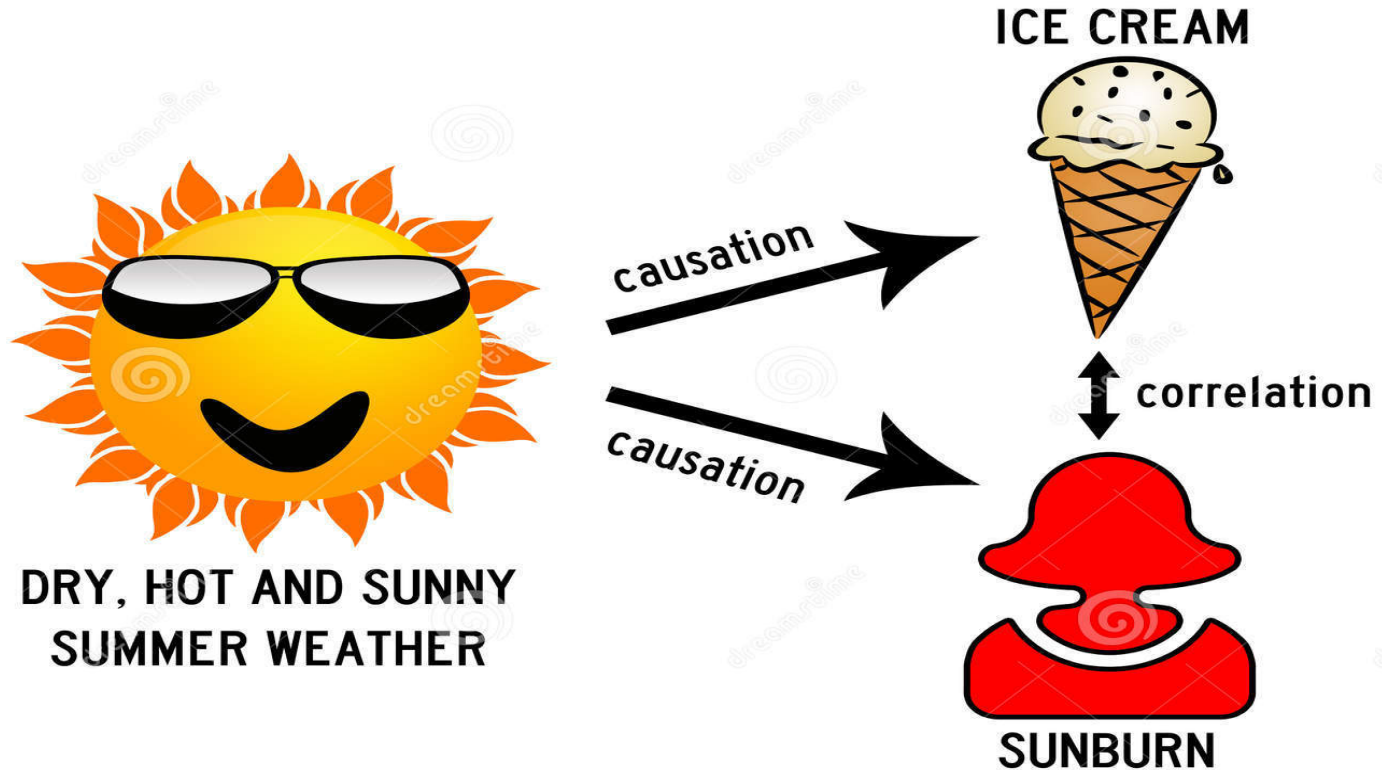


Cluster/group





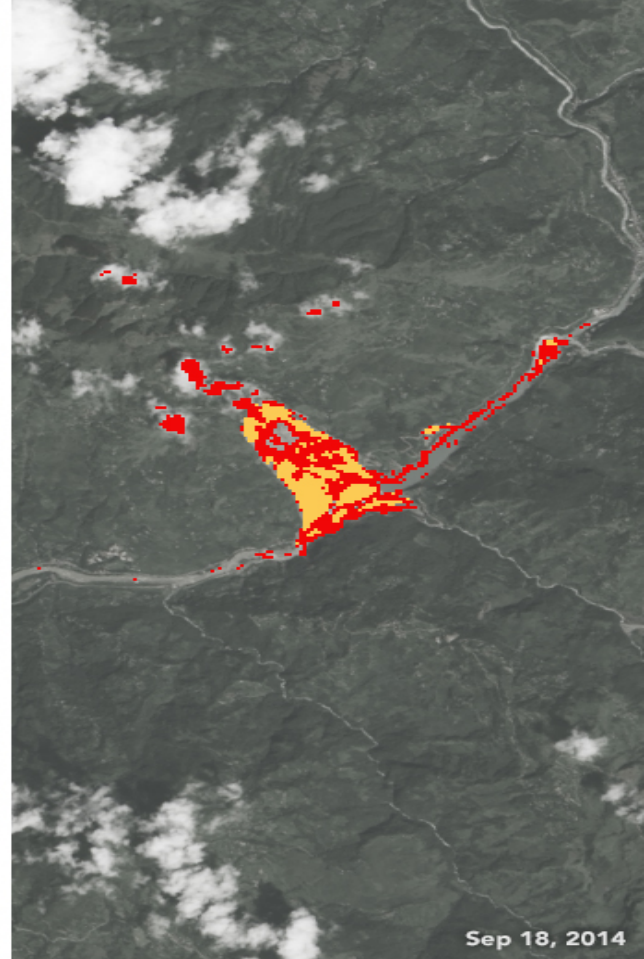
# Associations, A & B & C occur frequently





# Summarization





**Landslide Detection**

Possible

Probable

# Project



# Where do I get data?

- IVOA
  - VO
- VLA, Vizier, NED, ...
  - Extract all AGN from Veron paper within  $0 < \text{dec} < -90$
  - Plot the redshift distribution of the sources you have obtained

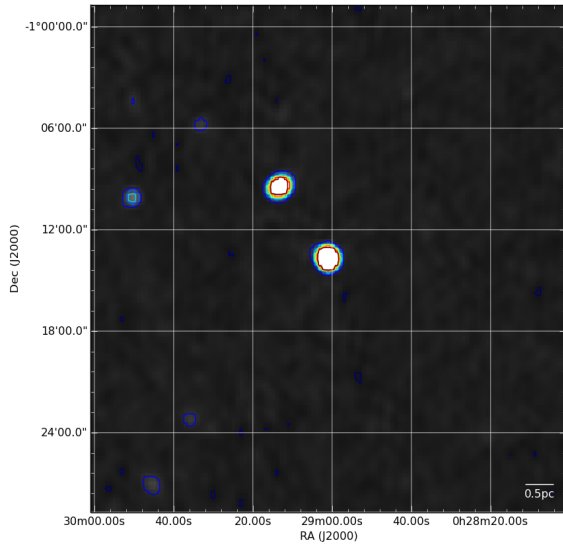
# Use Case



Morphology  
of  
Active Galactic Nuclei

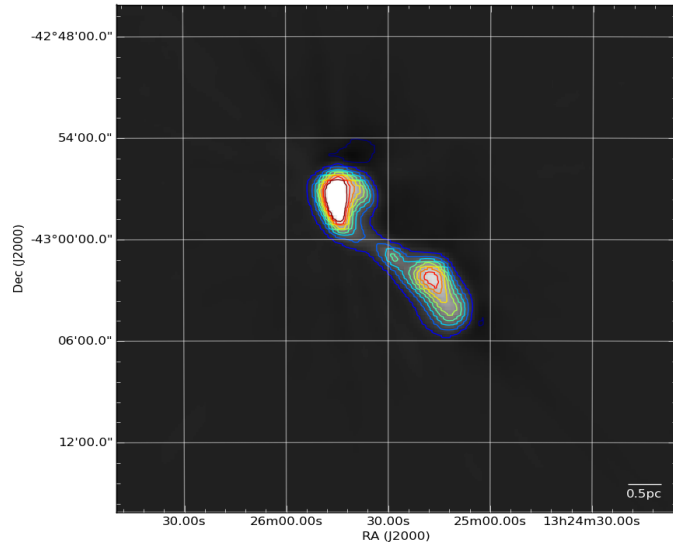
# Radio Sources

## Point Sources



PKS 0026-014

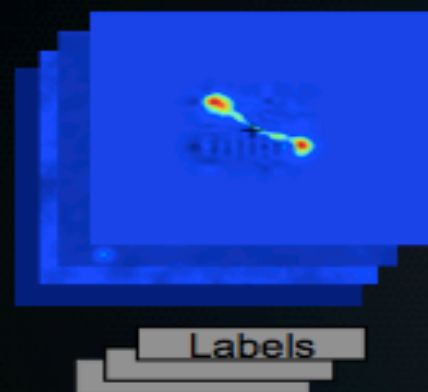
## Extended Sources



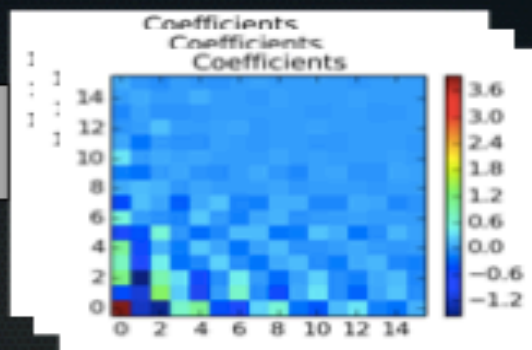
Centaurus A

# ML Classification

## Feature Extraction



Shapelets  
Decomposition



Feature  
Vectors



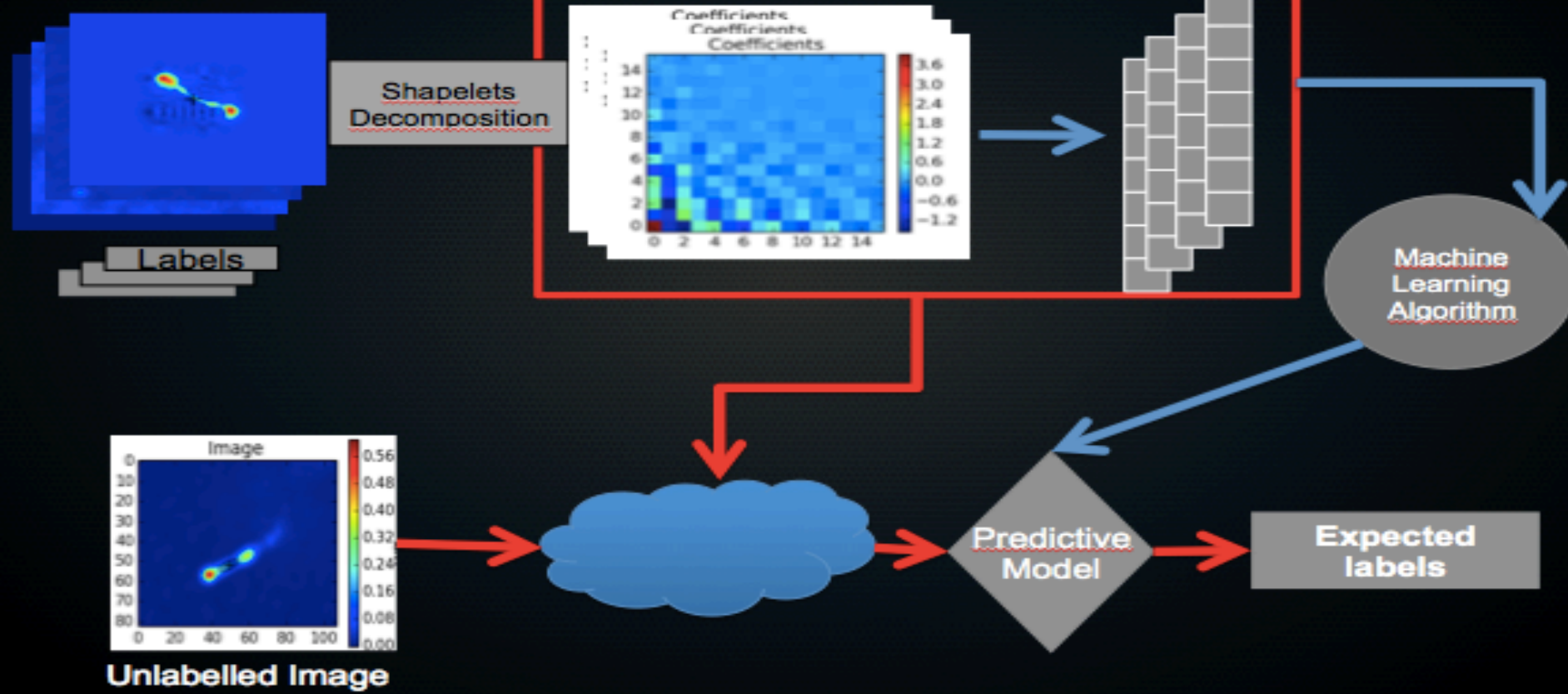
$$\mathbf{S}_{(i,\theta)} = \begin{bmatrix} f(0,0) & f(0,1) & f(0,2) & \dots & f(0,15) \\ f(1,0) & f(1,1) & f(1,2) & \dots & f(1,15) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ f(15,0) & f(15,1) & f(15,2) & \dots & f(15,15) \end{bmatrix}$$

$$\mathbf{x}_{(i,\theta)}^{\wedge} = \frac{\mathbf{X}_{(i,\theta)}}{|\mathbf{X}_{(i,\theta)}|} = \begin{bmatrix} g(0,0) \\ \vdots \\ g(0,15) \\ \vdots \\ \vdots \\ g(15,0) \\ \vdots \\ g(15,15) \end{bmatrix}$$



# ML Classification

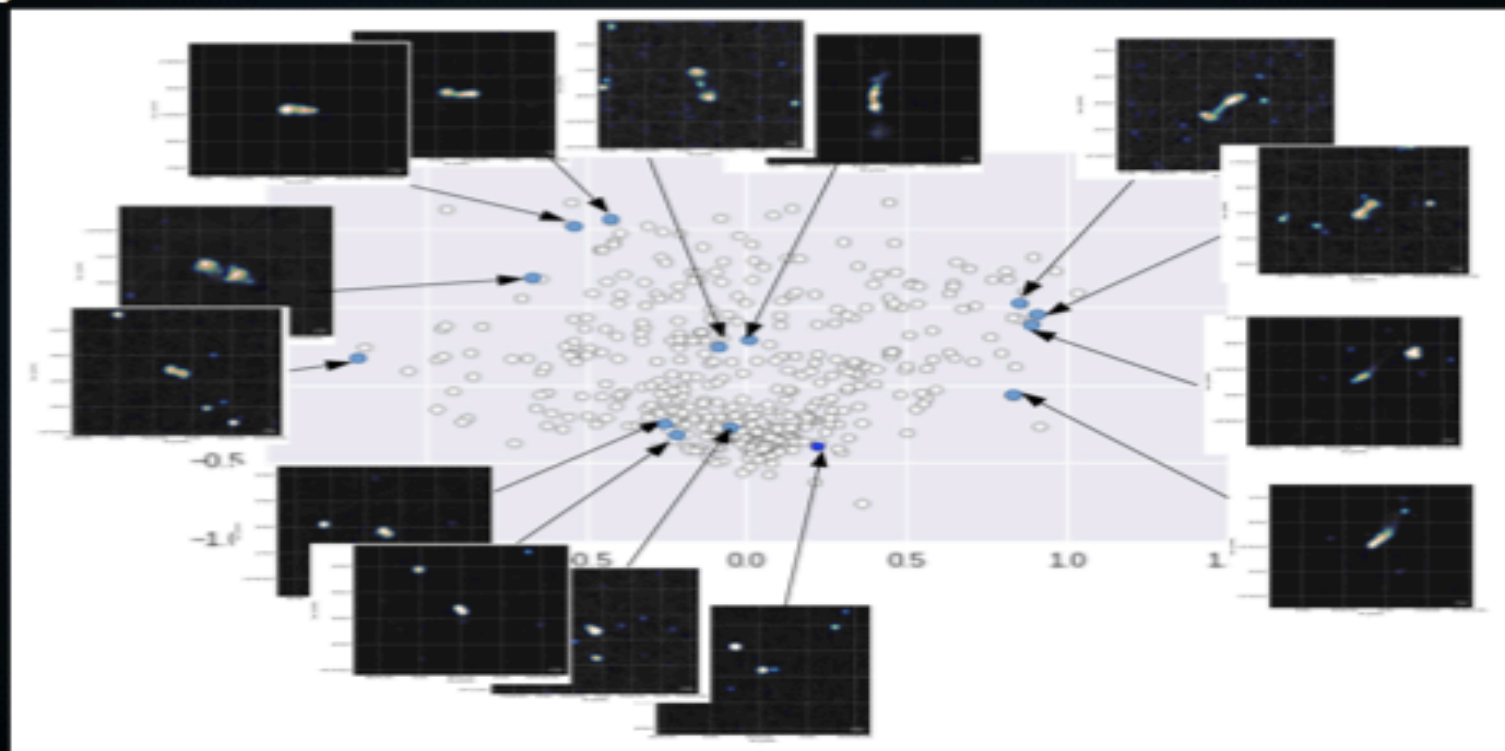
## Feature Extraction



# Feature Visualisation: isomap

Angular Variation

Triple Gaussian



Single Gaussian

Elliptical

# Use Case

- Morphology of AGN from FIRST
- Mine archive for list of AGN
- Images
  - Radio
  - Optical
  - Overlay
- Summary

Morphology of AGN

