# Data Mining Methods
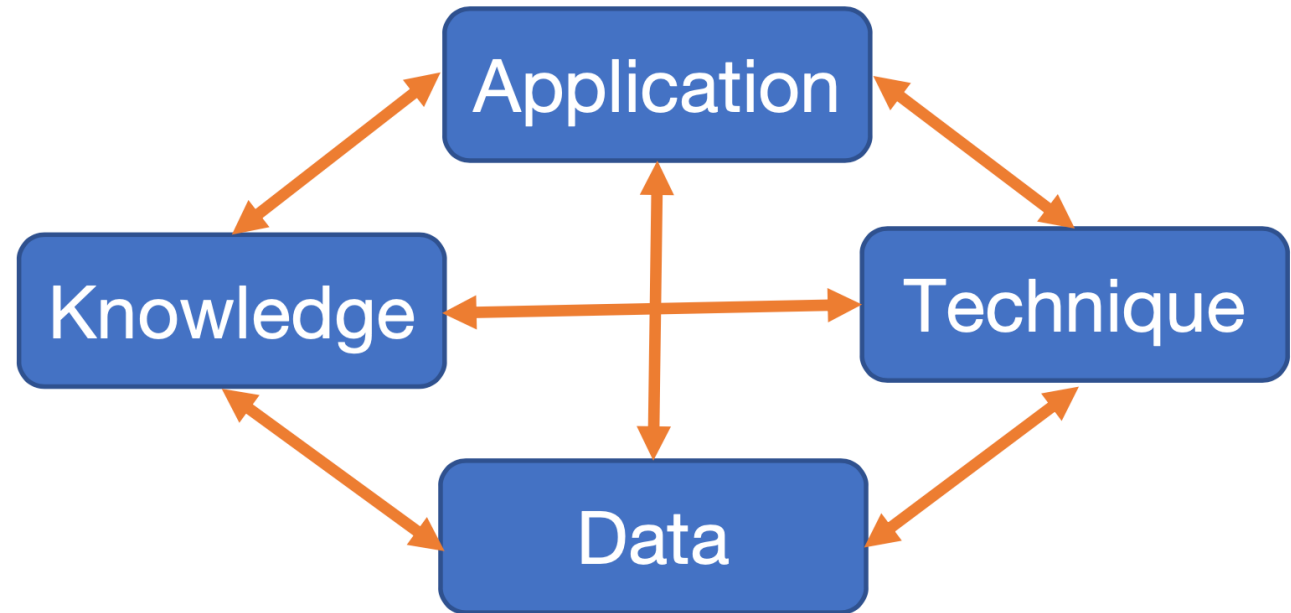## from data to insight

There is no such a thing as <span style="color:red">insight</span>
without a clear and concise question,
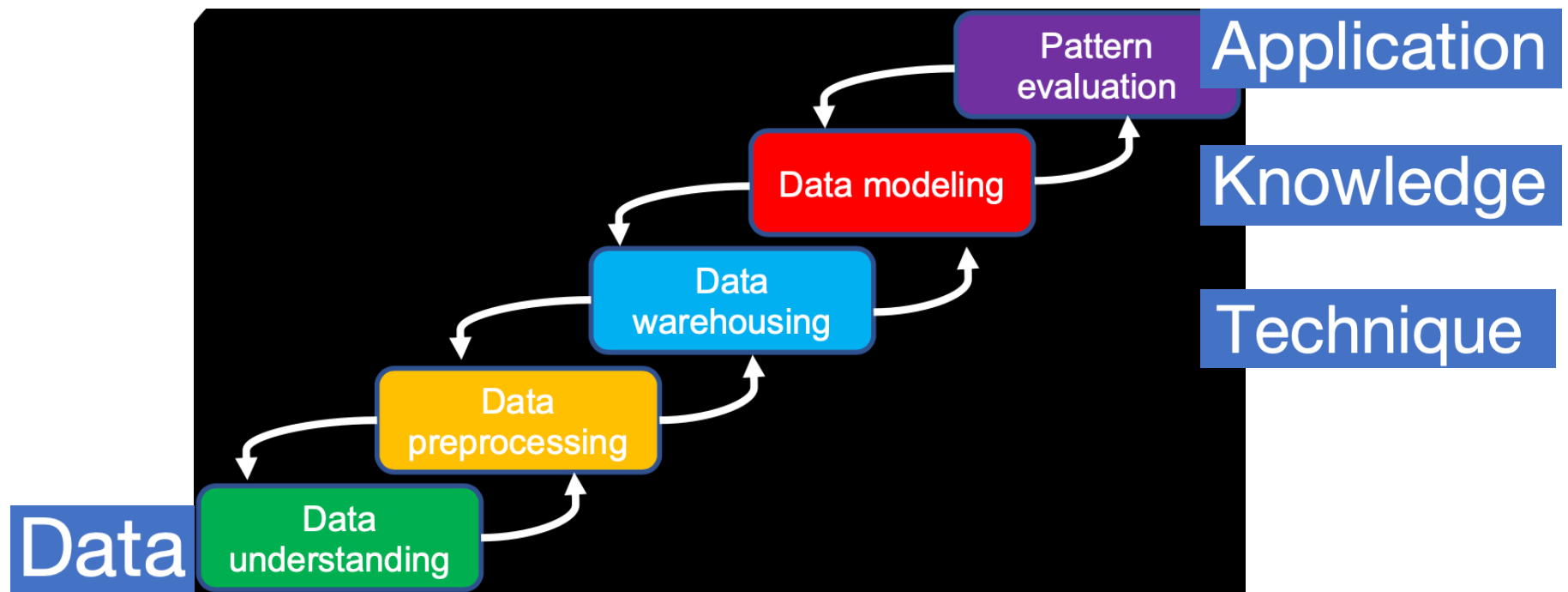as well as having a way to measure the success or failure
of the answer obtained

# Learning objective:

- o Identify the core functionalities of data modelling in the data mining pipeline.

- o Apply the Apriori algorithm for frequent itemset mining, among others

# Data Mining: Four Views

# Data Mining Pipeline

# Technique View

o   Frequent pattern analysis

o   Classification & prediction

o   Clustering

o   Anomaly detection

o   Trend and evolution analysis
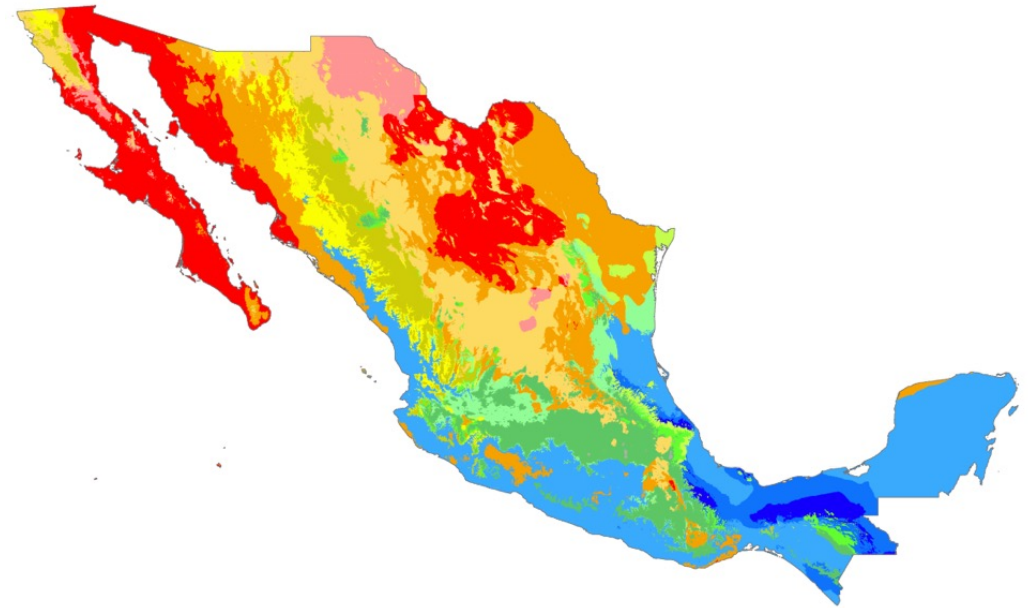
# Frequent Pattern Analysis

- o Frequent itemset

- o Frequent sequence

- o Frequent structure

- o Association rules

- o Correlation analysis

# Classification

o   Pre-defined classes

o   Need training data

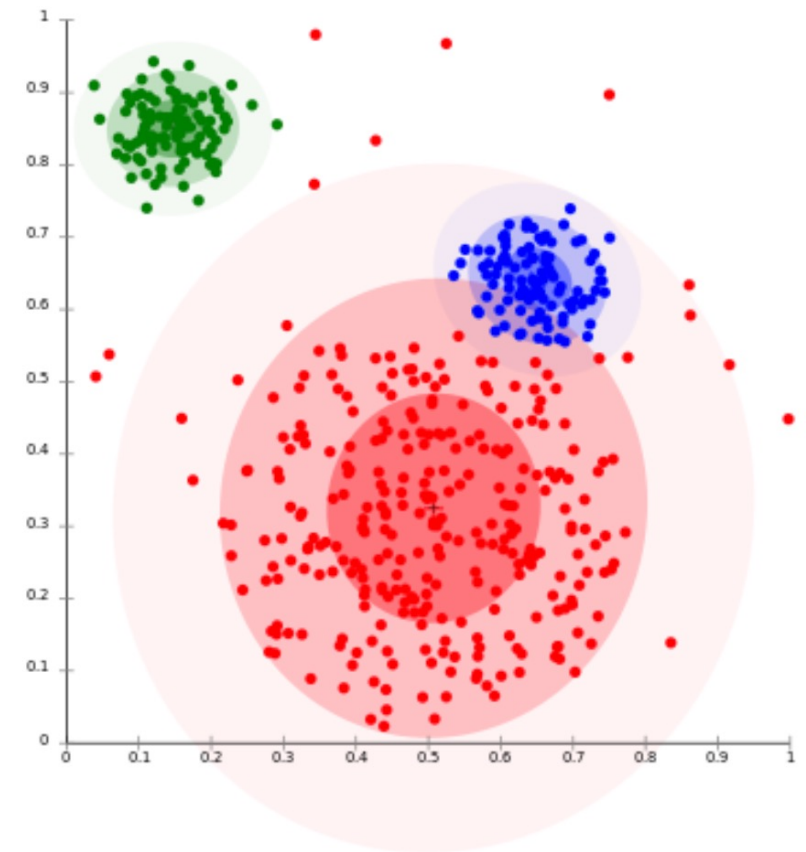o   Build model to distinguish classes

# Prediction

o   Numerical prediction (continuous value)

  • E.g., weather, stock price, traffic

# Clustering

o   No predefined classes

o   Intra-cluster similarity

o   Inter-cluster dissimilarity

# Anomaly Detection

o Anomaly/outlier - Differ from the "norm"

- E.g., error, noise, fraud, extreme events

**Daliana Liu** · Following

Data Scientist, "The Data Scientist Show" Podcast Host

17h · Edited · 🌐

Do NOT be a data scientist if you:

1. are easily frustrated when research doesn't yield results.
2. don't like to deal with vaguely defined problems.
3. only love the math and theory, but don't want to communicate with non-tech folks.

Note that I didn't say 'bad at math'. Everyone can learn tools to do math.

But to be a data scientist, you need the mindset to deal with ambiguity and uncertainty, so you can solve the business problem.

What are some other qualities you believe a data scientist must have?

*I share my about data science career here www.dalianaliu.com free to subscribe.

#datascience

# Trend and Evolution Analysis

o   Changes over time

✓  Overall trend

✓  Periodical patterns

✓  Anomalies

# Data Mining Methods

o  Frequent pattern analysis

o  Classification

o  Clustering

o  Outlier analysis

# Data science skills

Range from

o     programming to design

o     mathematics to storytelling

# The motivation – Data mining

o   Deriving valuable insights from data

    ✓   widely welcomed by businesses

# Typical questions :

o   What product will sell better in conjunction with another popular product?

o   How can customers be encouraged to spend a longer time in an online portal?

o   What advertisement should be placed on what site?

o   How to determine if a retail transaction is valid ?

# The skills

o   Both statistics and a strong business acumen

o   Foundations in computer science, mathematics, modelling and programming

o   Good communication skills & inquisitive mind

Inter alia

# Sexiest job of the 21st century

o data scientist

o it is hard to find people with the right skills to fill in these roles

o this has lead to branding

data scientists as Unicorns.

# Good data scientist

a linear combination of some of the following traits:

o    Curiosity

o    Grasp of machine learning

o    Data product building and management

o    Effective communication of data insights

o    Programming and data visualisation abilities

o    Knowledge of statistics and probability

o    Healthy skepticism, in the scientific tradition

# Four pillars

1. Identify who the main stakeholders and clarify the lines of reporting.

2. To be able to work independently and productively

3. Identify the data to tackle a problem

   o proper interpretation is not necessarily easy, and misrepresentation of the results can be very damaging.

4. Have the outcome always in mind

# Technologies

o   Data Framework - MapReduce, BigQuery, Hadoop, Spark

o   Streaming data collection - Kafka, Flume, Scribe

o   Job scheduling - Azkaban, Oozie

o   Big Data Query languages - Pig, Hive

o   Data stores - Voldemort, Cassandra, Neo4j, Hbase

# Open Source Tools

o   **Python**: Data manipulation, prototyping, scripting

o   **Apache Hadoop:** Framework for processing big data

o   Apache Mahout: Scalable machine-learning algorithms for Hadoop

o   Spark: Cluster-computing framework for data analytics

o   R Project for Statistical Computing: Data manipulation and graphing

o   Julia: High-performance technical computing

o   **GitHub**, Subversion: Software and model management tools

o   Ruby, Perl, OpenRefine: Prototyping and production scripting languages

# The steps

o   Question identification

o   Data acquisition

o   Data munging – wrangling -  data janitor

o   Model construction

o   Representation

o   Interaction