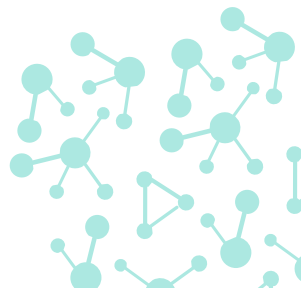# Data Science At AWS
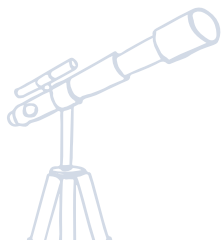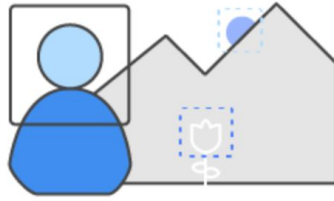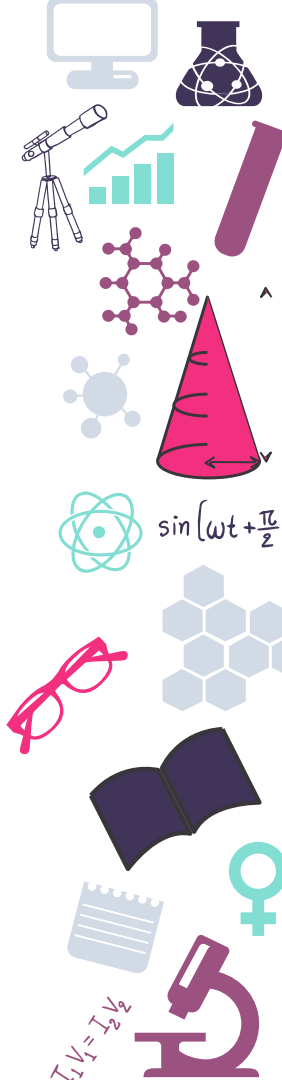
Rui LIU, Bruno BALDEZ CORREA

# 1

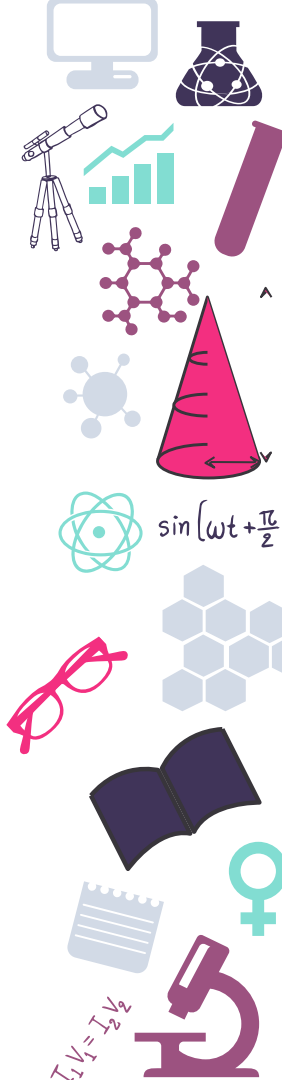## Data Science Services at AWS

# 2 Amazon Redshift

Amazon Redshift is a fast, fully managed, petabyte-scale cloud-base data warehouse solution offered by Amazon Web Services that provides simple and cost-effective functionalities to analyze all your data using standard SQL and BI techniques.

# Characteristics

- ⊘ Fast
- ⊘ Inexpensive
- ⊘ Extensible
- ⊘ Simple
- ⊘ Scalable
- ⊘ Secure
- ⊘ Compatible

# Architecture

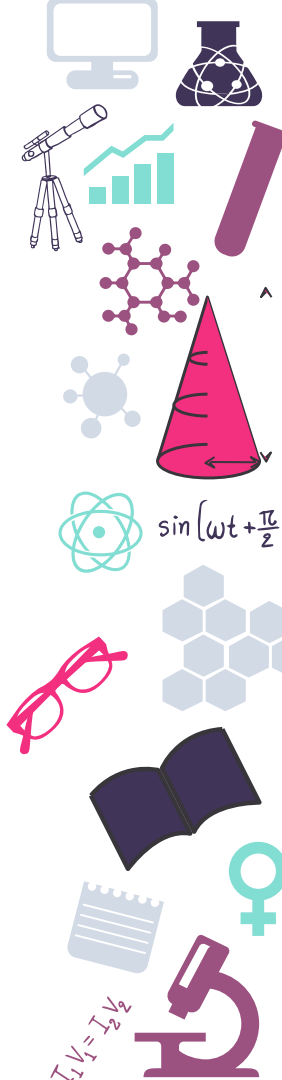## Clusters

- ✓ 1 or more compute nodes.

## Leader node

- ✓ compiles code
- ✓ distributes it

## Compute nodes
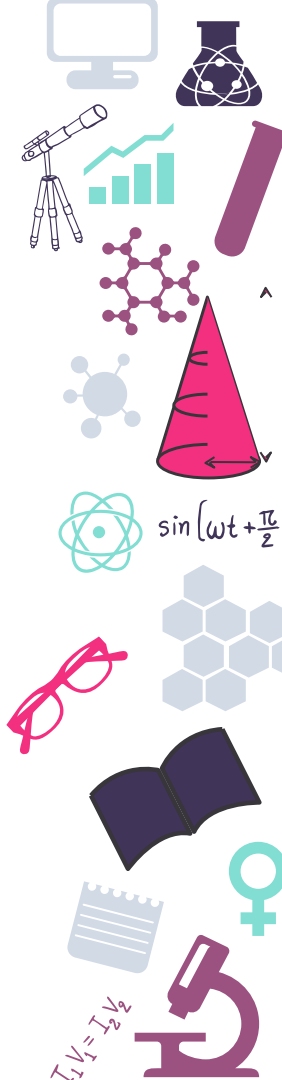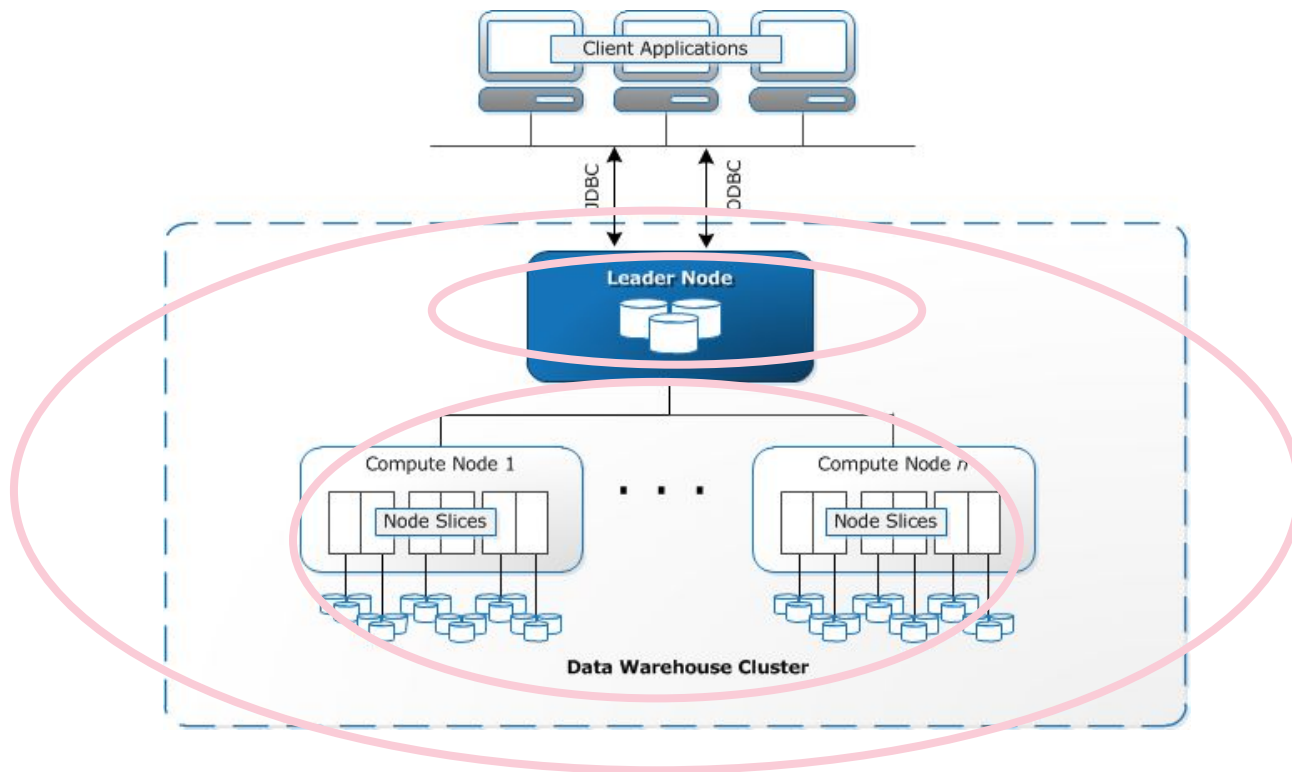
- ✓ compiles code
- ✓ execute the compiled code
- ✓ send intermediate results
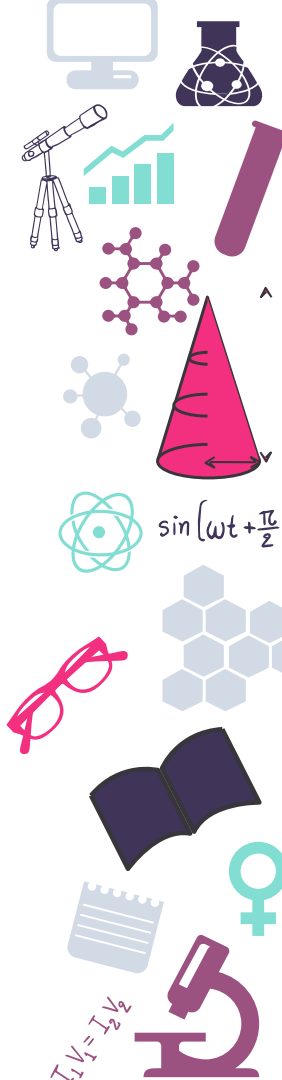- ✓ **dense storage nodes** or **dense compute nodes**

# Architecture

# Data Distribution

- [x] ALL
- [x] EVEN
- [x] KEY


- [x] SortKey
- [x] Primary and Foreign keys

# Getting Started with Amazon Redshift

**Step 1:** Set Up Prerequisites
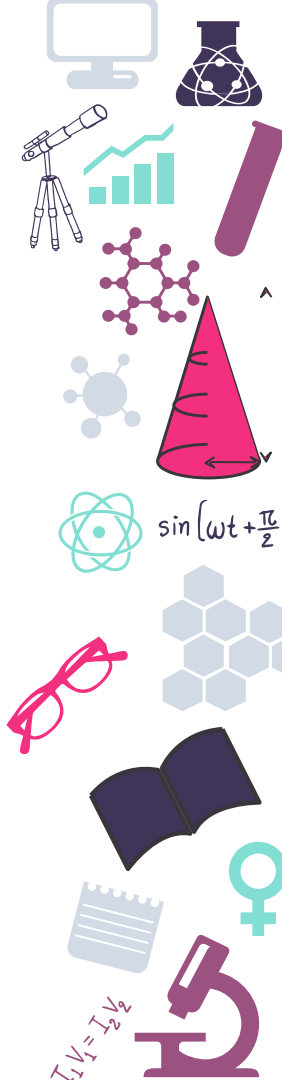**Step 2:** Create an IAM Role
**Step 3:** Launch a Sample Amazon Redshift Cluster
**Step 4:** Authorize Access to the Cluster
**Step 5:** Connect to the Sample Cluster
**Step 6:** Load Sample Data from Amazon S3
**Step 7:** Find Additional Resources and Reset Your Environment

# Redshift Use Scenario

# Redshift Use Scenario

# Redshift Use Scenario

# 3

# Amazon Machine Learning

Amazon Machine Learning makes it easy for developers to build machine learning model without learning complex algorithm or hiring experts.

# Ideal Usage Pattern

### Hard to code rules

- Rules are not explicit
- Number of factors are huge

### Hard to scale

- Large number of tasks
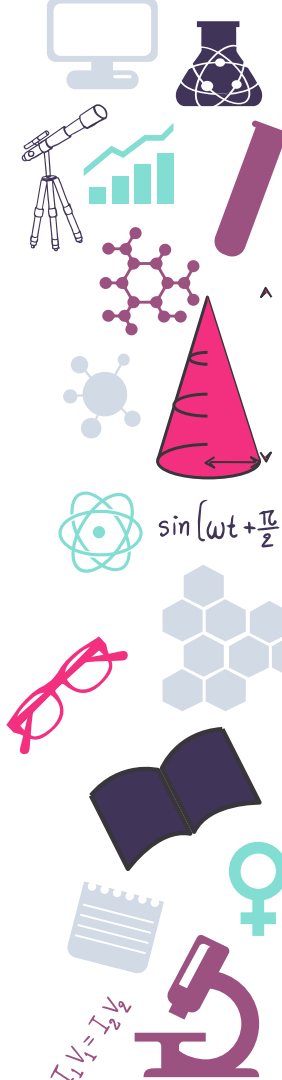- Impossible to classify tasks manually

# Datasources

Datasource is an object used by Amazon Machine Learning as train data, evaluation data and validation data
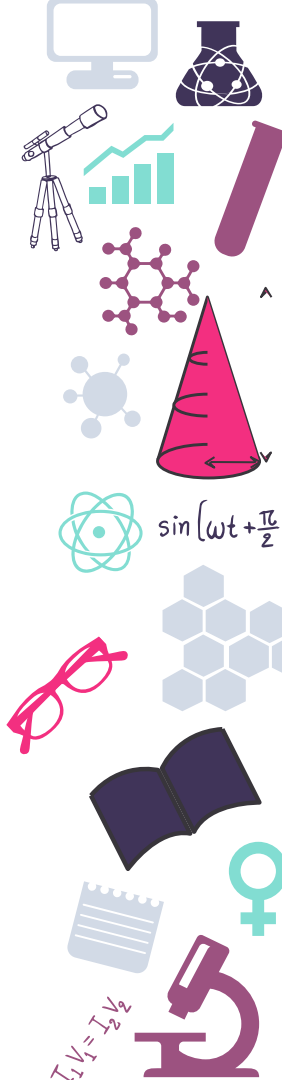
- ✓  Data should be well-formatted
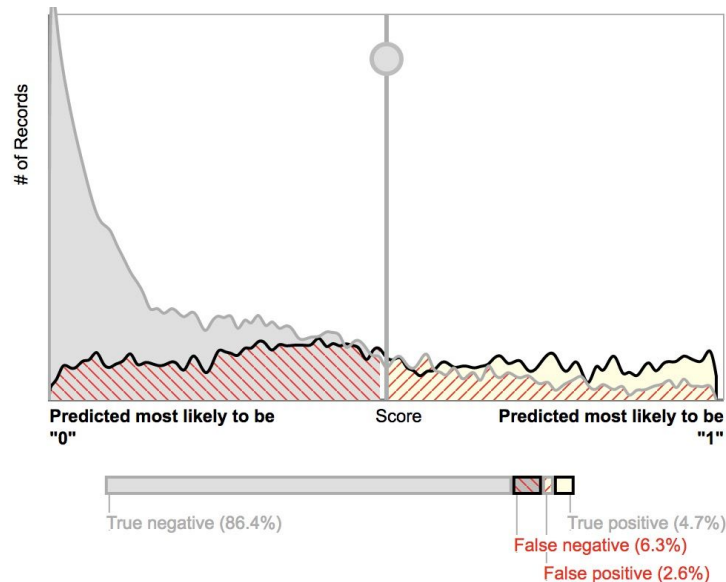- ✓  Datasources should contain one column as target

# Train ML Models

Amazon ML applies machine learning algorithms automatically

- ✓ Binary Classification
- ✓ Multiclass Classification
- ✓ Regression Model

# Evaluate ML Models
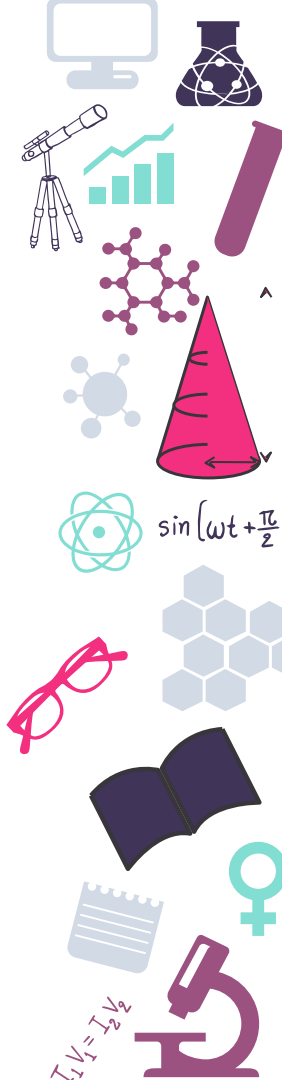
# Generate and Interpret Prediction

- ✓ Batch Prediction
- ✓ Real-Time Prediction

# Cost Model

**Size of Model**

**Number of Predictions**

# Advantages & Disadvantages

## Advantages

- Automatic
- Fast and easy

## Disadvantages

- Black box
- Supervised model only

# Thanks!

**Any questions?**