**Data Science with Python**

Machine Learning with Scikit–Learn

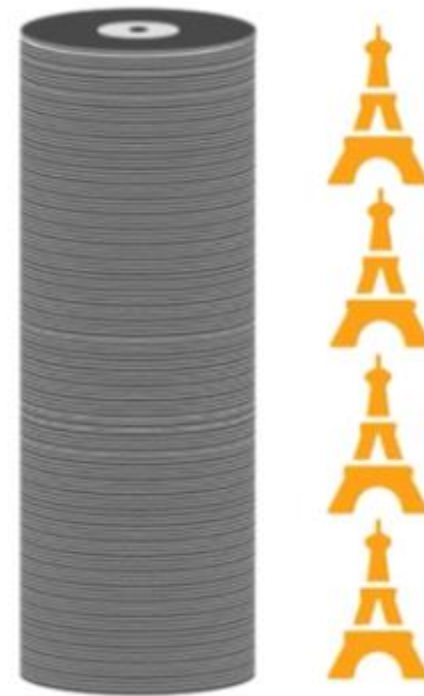# Learning Objectives

By the end of this lesson, you will be able to:

- Define machine learning

- Explain the machine learning approach

- List relevant terminologies that help you understand a dataset

- Discuss features of supervised and unsupervised learning models

- Explain algorithms, such as regression, classification, clustering, and dimensionality reduction

# Introduction to Machine Learning

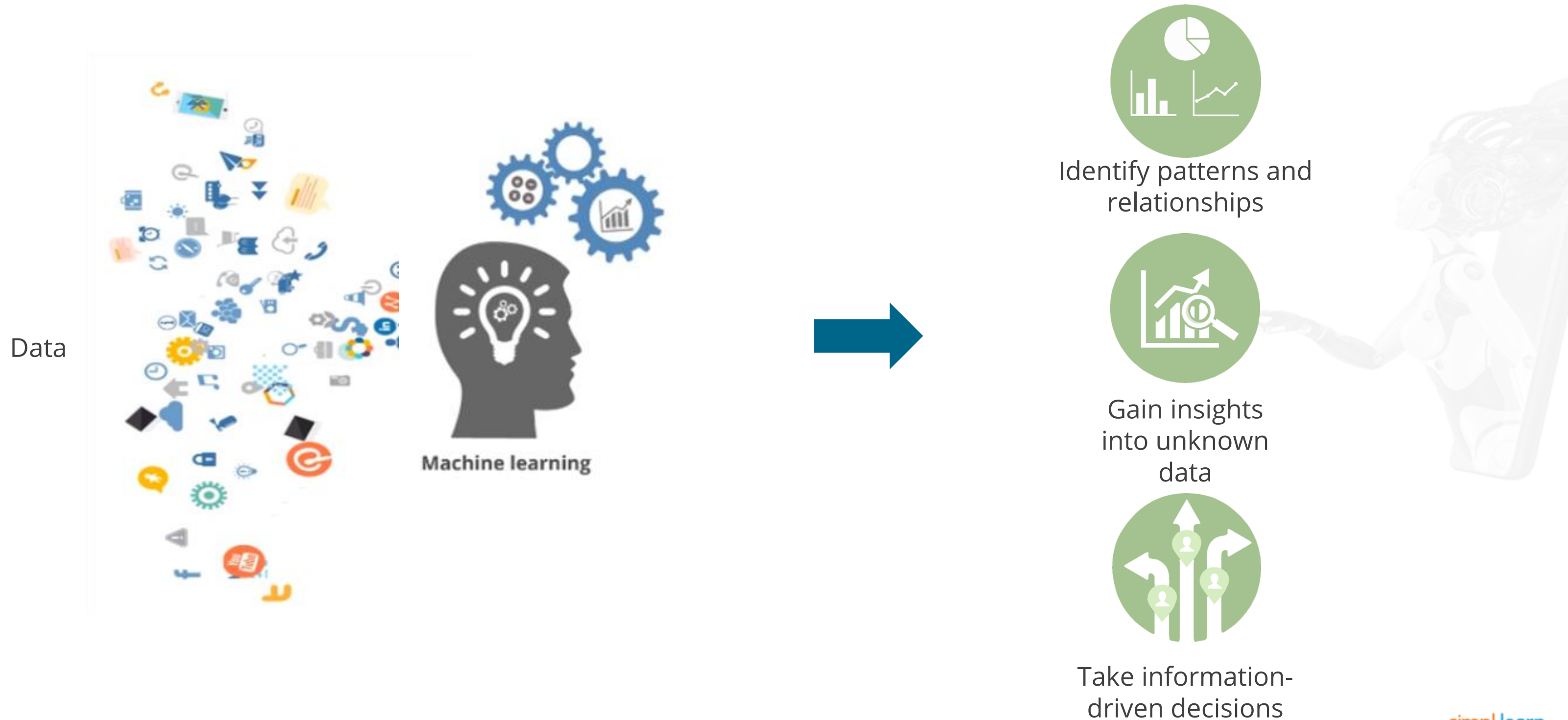simpli learn

# Why Machine Learning

If we stored the data generated in a day on Blu-ray disks and stacked them up, it would be equal to the height of four Eiffel towers. Machine learning helps analyze this data easily and quickly.



Machine learning

# Purpose of Machine Learning

Machine learning is a great tool to analyze data, find hidden data patterns and relationships, and extract information to enable information-driven decisions and provide insights.
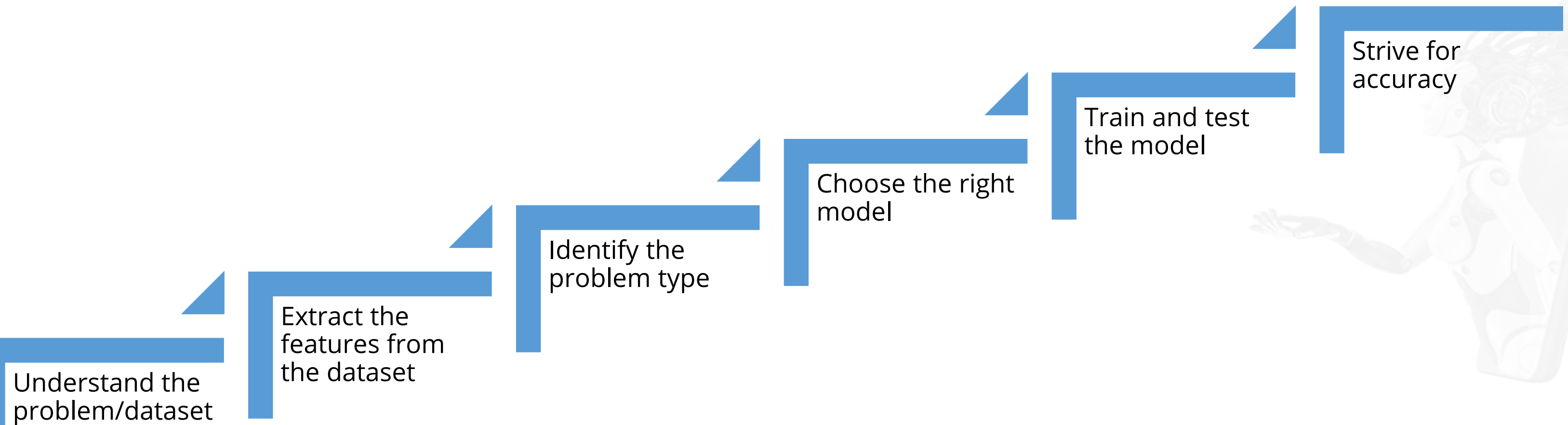
Data

Machine learning

Identify patterns and relationships

Gain insights into unknown data

Take information-driven decisions

# Machine Learning Terminology

These are some machine learning terminologies that you will come across in this lesson:

Inputs

Attributes

Features

Label

Records

Response

Observations

Outcome

Examples

Target

Samples

# Machine Learning Approach

# Machine Learning Approach

The machine learning approach starts with either a problem that you need to solve or a given dataset that you need to analyze.

Strive for accuracy

Train and test the model

Choose the right model

Identify the problem type

Extract the features from the dataset

Understand the problem/dataset

# Steps 1 and 2: Understand the Dataset and Extract Its Features

Let us look at a dataset and understand its features in terms of machine learning.

Features (attributes)

Response (label)

| Education (Yrs.) | Professional Training (Yes/No) | Hourly Rate (USD) |
|---|---|---|
| 16 | 1 | 90 |
| 15 | 0 | 65 |
| 12 | 1 | 70 |
| 18 | 1 | 130 |
| 16 | 0 | 110 |
| 16 | 1 | 100 |
| 15 | 1 | 105 |
| 31 | 0 | 70 |

Observations (records)

Predictors

# Steps 3 and 4: Identify the Problem Type and Learning Model

Machine learning can either be supervised or unsupervised. The problem type should be selected based on the type of learning model.

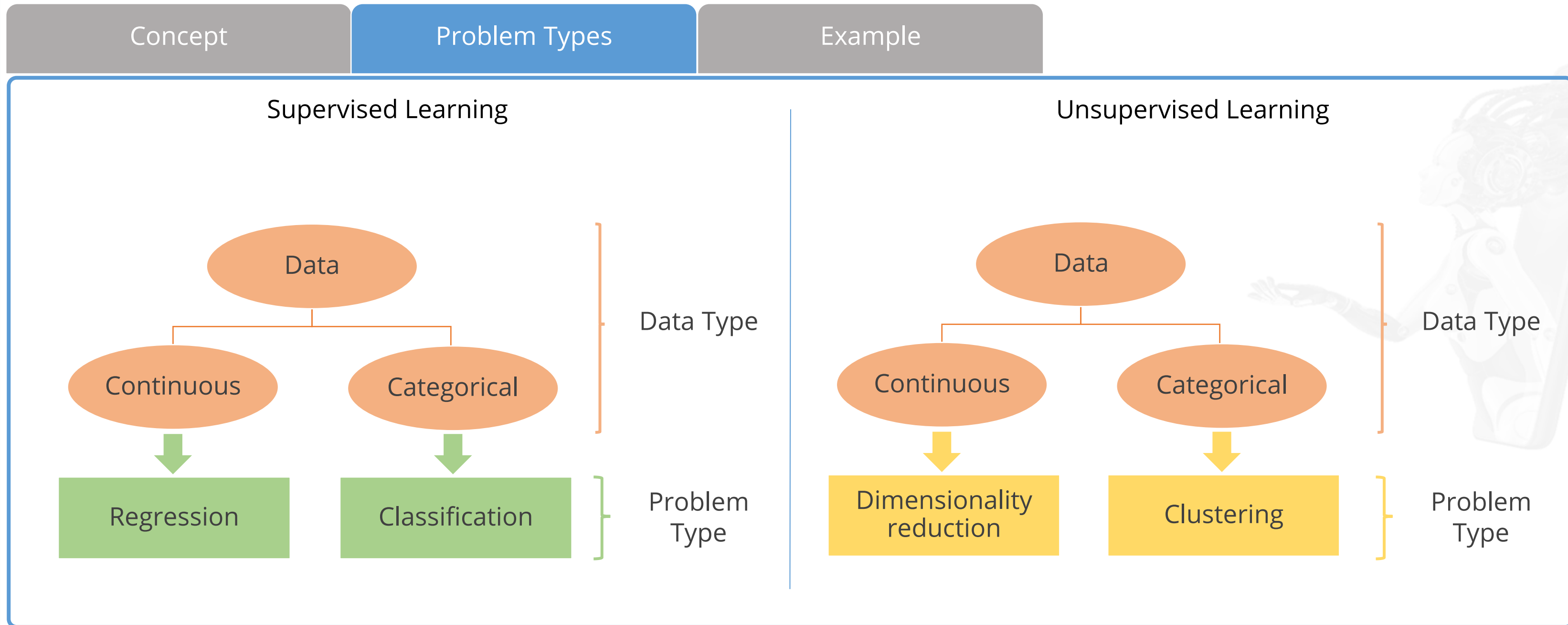| Concept | Problem Types | Example |
| --- | --- | --- |

### Supervised Learning

- In supervised learning, the dataset used to train a model should have observations, features, and responses. The model is trained to predict the "right" response for a given set of data points.

- Supervised learning models are used to predict an outcome.

- The goal of this model is to "generalize" a dataset so that the "general rule" can be applied to new data as well.

### Unsupervised Learning

- In unsupervised learning, the response or the outcome of the data is not known.

- Unsupervised learning models are used to identify and visualize patterns in data by grouping similar types of data.

- The goal of this model is to "represent" data in a way that meaningful information can be extracted.

# Steps 3 and 4: Identify the Problem Type and Learning Model

Data can either be continuous or categorical. Based on whether it is supervised or unsupervised learning, the problem type will differ.

| Concept | Problem Types | Example |
|---|---|---|

## Supervised Learning

```
              Data
          /          \
   Continuous      Categorical        } Data Type
       |                |
       ↓                ↓
  Regression      Classification      } Problem Type
```

## Unsupervised Learning

```
              Data
          /          \
   Continuous      Categorical        } Data Type
       |                |
       ↓                ↓
  Dimensionality    Clustering        } Problem Type
   reduction
```

# Steps 3 and 4: Identify the Problem Type and Learning Model

Some examples of supervised and unsupervised learning models are:

| Concept | Problem Types | Example |
|---------|---------------|---------|

### Supervised Learning
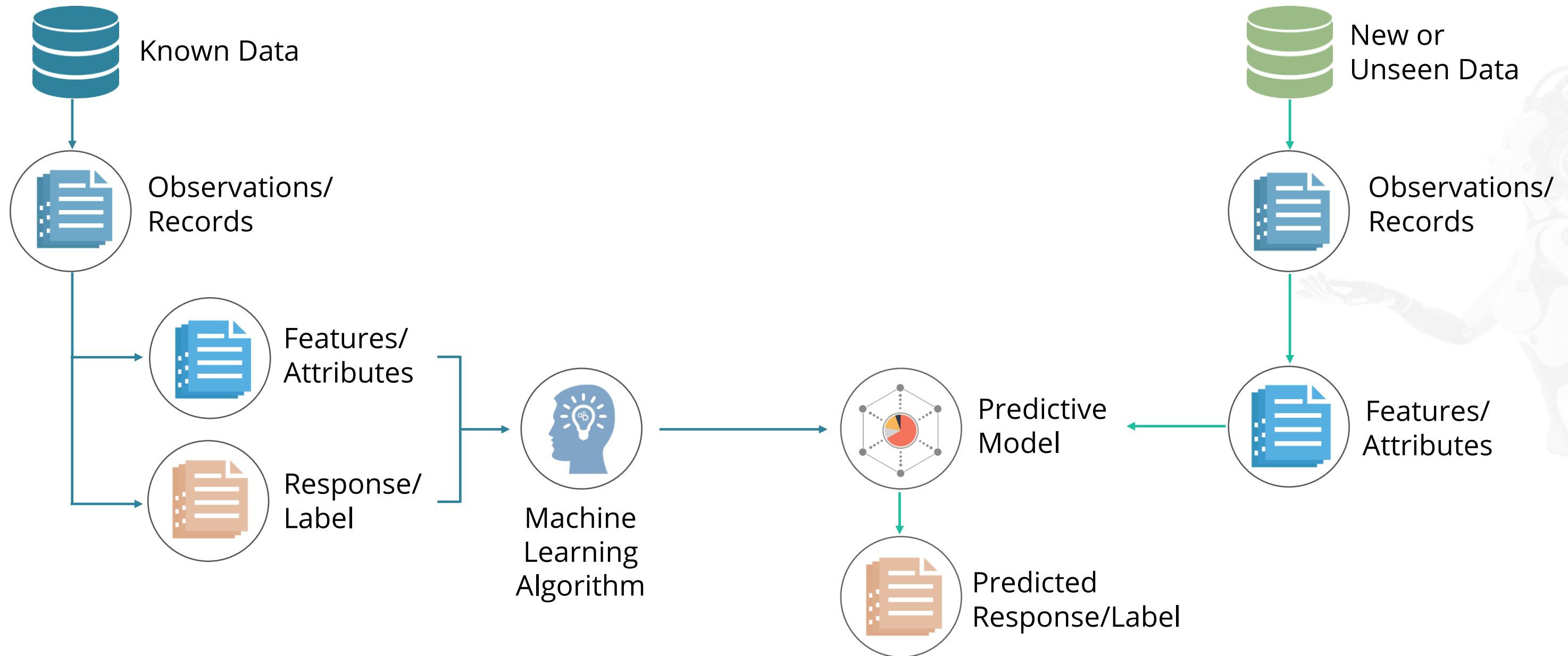


Categories of news based on the topics

### Unsupervised Learning
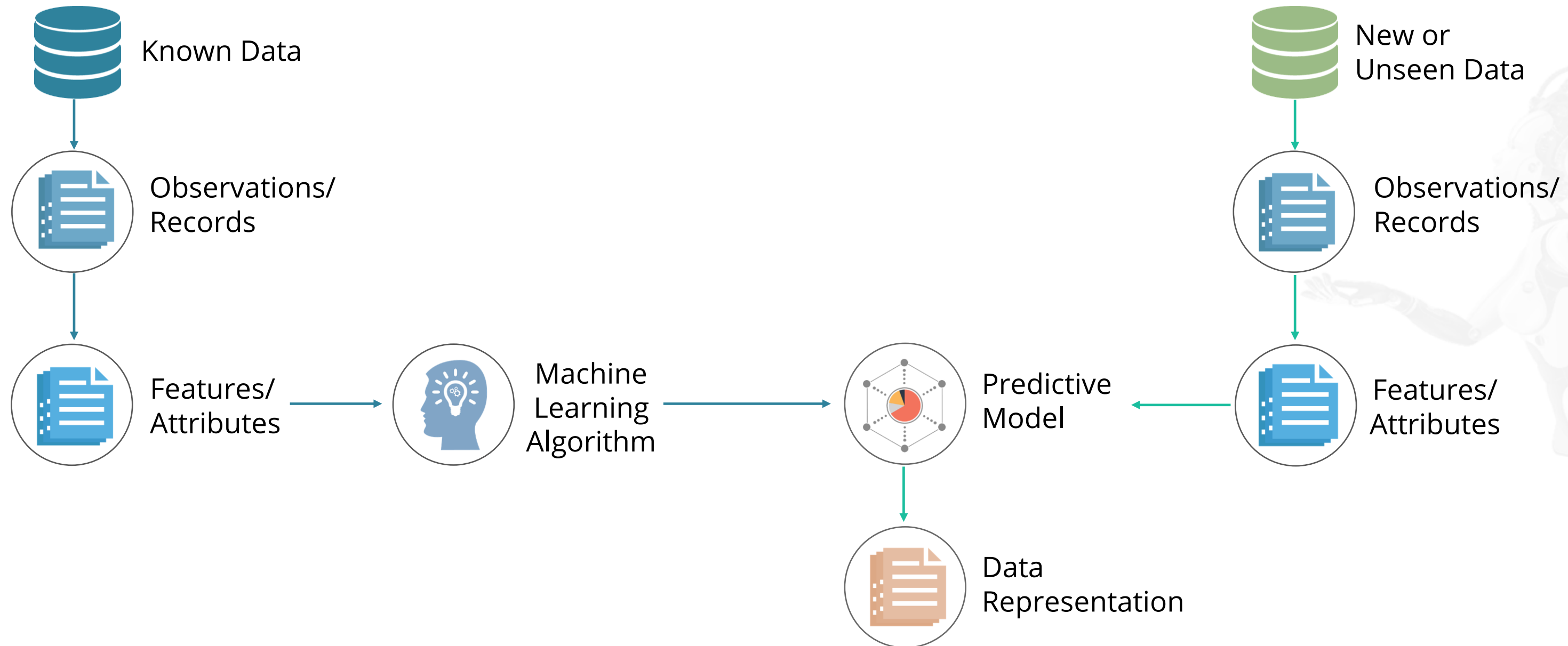


Grouping of similar stories on different news networks

# Working of Supervised Learning Model

In supervised learning, a known dataset with observations, features, and response is used to create and train a machine learning algorithm. A predictive model, built on top of this algorithm, is then used to predict the response for a new dataset that has the same features.

Known Data

Observations/
Records

Features/
Attributes

Response/
Label

Machine
Learning
Algorithm

Predictive
Model

Predicted
Response/Label

New or
Unseen Data

Observations/
Records

Features/
Attributes

# Working of Unsupervised Learning Model

In unsupervised learning, a known dataset has a set of observations with features. But the response is not known. The predictive model uses these features to identify how to classify and represent the data points of new or unseen data.

Known Data

Observations/ Records

Features/ Attributes

Machine Learning Algorithm

Predictive Model

Data Representation

New or Unseen Data

Observations/ Records

Features/ Attributes

# Steps 5 and 6: Train, Test, and Optimize the Model

To train supervised learning models, data analysts usually divide a known dataset into training and testing sets.

## Supervised Learning

Known Data

Observations/ Records

Features/ Attributes

Response/ Label

## Unsupervised Learning

Known Data

Observations/ Records

Features/ Attributes

# Steps 5 and 6: Train, Test, and Optimize the Model



Known Data

Train
(60%-80%)

Test
(20%-40%)

Observations/
Records

Features/
Attributes

Response
/ Label

Machine
Learning
Algorithm

# Steps 5 and 6: Train, Test, and Optimize the Model

Let us look at an example to see how the split approach works.

Model Training

Observation

Response

Train set

Test set

Train set

Test set

| ID | Education | Professional training | Hourly rate |
|----|-----------|-----------------------|-------------|
| 10 | 16 | 1 | 90 |
| 45 | 15 | 0 | 65 |
| 83 | 12 | 1 | 70 |
| 45 | 18 | 1 | 130 |
| 54 | 16 | 0 | 110 |
| 67 | 16 | 1 | 100 |
| 71 | 15 | 1 | 105 |
| 31 | 15 | 0 | 70 |

simplilearn

# Supervised Learning Model Considerations

Some considerations of supervised and unsupervised learning models are shown here.

# Scikit-Learn

# Scikit-Learn

Scikit is a powerful and modern machine learning Python library for fully and semi-automated data analysis and information extraction.

Efficient tools to identify and organize problems (Supervised/ Unsupervised)

Free and open datasets

Rich set of libraries for learning and predicting

Model support for every problem type

Model persistence

open source initiative ®

Open source community and vendor support

# Scikit-Learn: Problem-Solution Approach

Scikit-learn helps Data Scientists organize their work through its problem-solution approach.

| Model Selection | Estimator Object | Model Training | Predictions | Model Tuning | Accuracy |

# Scikit-Learn: Problem-Solution Considerations

While working with a Scikit-Learn dataset or loading your own data to Scikit -Learn, always consider these points:

✔ Create separate objects for feature and response

✔ Ensure that features and response have only numeric values

✔ Features and response should be in the form of a NumPy ndarray

✔ Since features and response would be in the form of arrays, they would have shapes and sizes

✔ Features are always mapped as $x$, and response is mapped as $y$

**Supervised Learning Models: Linear Regression**

# Supervised Learning Models: Linear Regression

Linear regression is a supervised learning model used to analyze continuous data.

It is easy to use as the model does not require a lot of tuning

It is the most basic and widely used technique to predict a value of an attribute

It runs very fast, which makes it more time-efficient

# Supervised Learning Models: Linear Regression

The linear regression equation is based on the formula for a simple linear equation.

$$y = mx + c$$

Simple linear equation

$$y = \beta_0 + \beta_1 x$$

Linear regression equation

Response

Input features

Intercept

Coefficient of x

# Supervised Learning Models: Linear Regression

Linear regression is the most basic technique to predict a value of an attribute.

Data point

Residual

Least square line

$y$ (response)

dy

dx

$\beta_1 = \dfrac{\mathrm{d}y}{\mathrm{d}x}$

Slope/ gradient

$(0, y)$

$\beta_0$

$x$ (predictor variable)

Residual

$$y = \beta_0 + \beta_1 x + u$$

Actual value

Predicted value

**!** The attributes are usually fitted using the "least square" approach.

# Supervised Learning Models: Linear Regression

Smaller the value of SSR or SSE, the more accurate the prediction will be, which would make the model the best fit.

Data point

Least square line

$$y = \beta_0 + \beta_1 x + u$$

$$SSR = \sum (\widehat{y_i} - \bar{y})^2$$

Regression of sum of squares

$$SSE = \sum (y_i - \widehat{y_i})^2$$

Error of sum of squares

SSE

SSR

$$\bar{y} = \bar{\beta}_0 + \bar{\beta}_1 x$$

$y$ (response)

$\bar{y}$

$(0, y)$

$\beta_0$

$x$ (predictor variable)

**!** The attributes are usually fitted using the "least square" approach.

# Supervised Learning Models: Linear Regression

Let us see how linear regression works in Scikit-Learn.

Normalizes the regression variable before performing the regression operation

Calculates the intercept for this model

Class

sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=1)

Copies the regression variable

Number of jobs to use for the computation

# Loading a Dataset

**Problem Statement:** Demonstrate how to load a built-in SciKit-learn dataset.

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

# Linear Regression Model

**Problem Statement:** Demonstrate how to create and train a linear regression model

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.
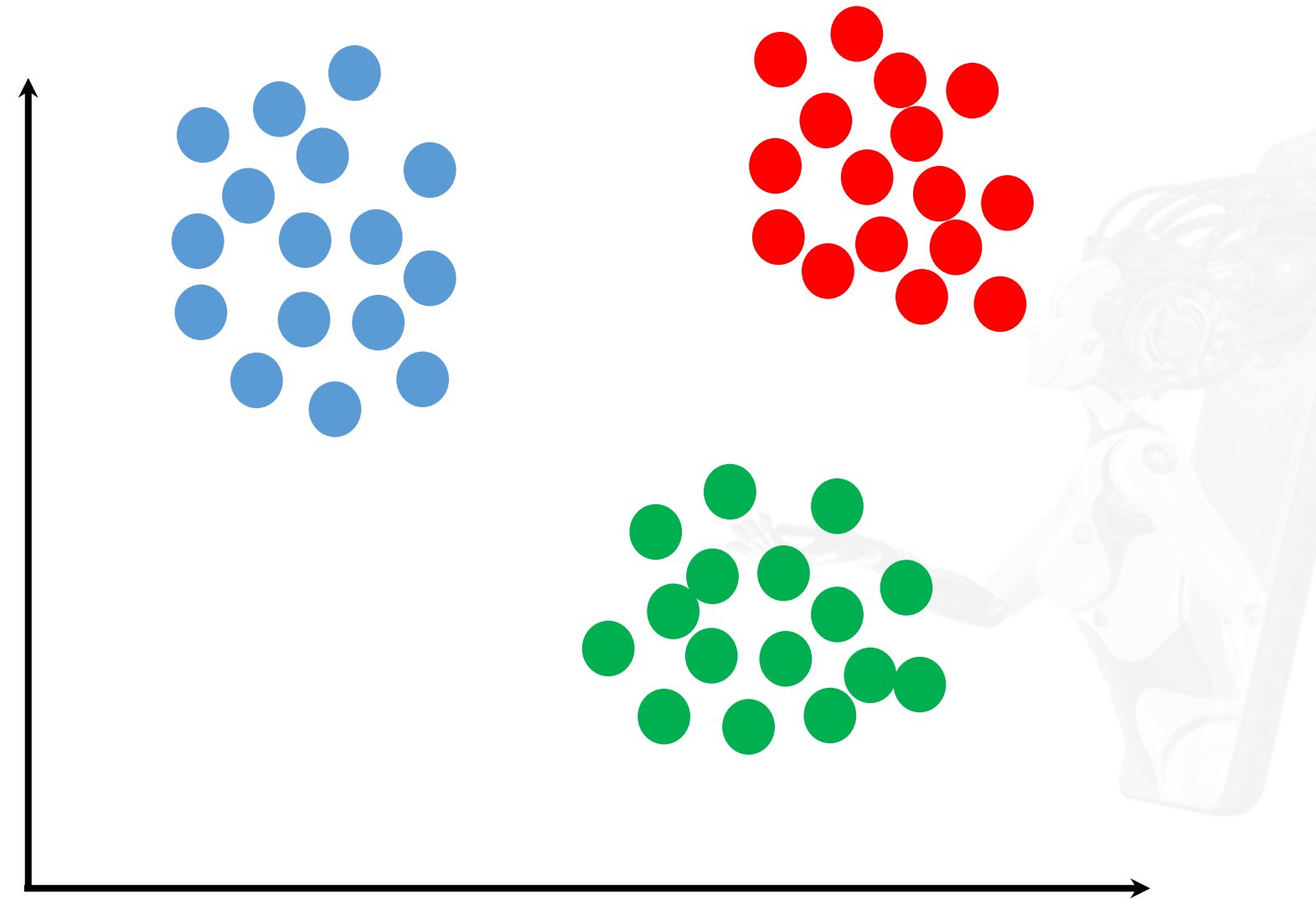
ASSISTED PRACTICE

simplilearn

Supervised Learning Models: Logistic Regression

Logistic regression is a generalization of the linear regression model used for classification problems.

$$\pi = \Pr(y = 1|x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

Probability of $y = 1$, given $x$

Change in the log-odds for a unit change in x

The purpose of K-NN is to predict the class for each observation.

# Supervised Learning Models: Logistic Regression

Logistic regression is a generalization of the linear regression model used for classification problems.

$$\text{Odds} = \frac{\pi}{1 - \pi}$$

Probability

$$\log\left(\frac{\pi}{1 - \pi}\right) = \log\left(e^{\beta_0 + \beta_1 x}\right) = \beta_0 + \beta_1 x$$

Logarithm of odds

Linear regression

# Supervised Learning Models: Logistic Regression

Specifies the norm used
in penalization

Inverse of
regularization

Calculates
the
intercept

Class

Implemented only
for L2 penalty

class sklearn.linear_model.LogisticRegression(penalty='l2', dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='liblinear', max_iter=100, multi_class='ovr', verbose=0, warm_start=False, n_jobs=1)

Seed or the random
state instance

Algorithm to use in the
optimization problem

If true, reuse the
solution of the
previous call

Can be ovr (binary)
or multinomial

Number of
jobs in parallel
computation

Supervised Learning Models: : K-Nearest Neighbors

# Supervised Learning Models: K-Nearest Neighbors (K-NN)

K-Nearest Neighbors, or K-NN, is one of the simplest machine learning algorithms used for both classification and regression problem types.

Features
(Attributes)

| Education (Yrs.) | Professional Training (Yes/No) | Hourly Rate (USD) |
|---|---|---|
| 16 | 1 | 90 |
| 15 | 0 | 65 |
| 12 | 1 | 70 |
| 18 | 1 | 130 |
| 16 | 0 | 110 |
| 16 | 1 | 100 |
| 15 | 1 | 105 |
| 31 | 0 | 70 |

# Supervised Learning Models: K-Nearest Neighbors



K=3

K=6

If you are using this method for binary classification, choose an odd number for k to avoid the case of a **tied** distance between two classes.

# Supervised Learning Models: K-Nearest Neighbors

K-means finds the best centroids by alternatively assigning random centroids to a dataset and selecting mean data points from the resulting clusters to form new centroids. It continues this process iteratively until the model is optimized.

Features (Attributes)

Response (label)

| Education (Yrs.) | Professional Training (Yes/No) | Hourly Rate (USD) |
|---|---|---|
| 16 | 1 | 90 |
| 15 | 0 | 65 |
| 12 | 1 | 70 |
| 18 | 1 | 130 |
| 16 | 0 | 110 |
| 16 | 1 | 100 |
| 15 | 1 | 105 |
| 31 | 0 | 70 |

# K-NN and Logistic Regression Models

**Problem Statement:** Demonstrate the use of K-NN and logistic regression models

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

ASSISTED PRACTICE

# Unsupervised Learning Models: Clustering

# Unsupervised Learning Models: Clustering

A cluster is a group of similar data points.

Clustering is used to:

- Extract the structure of the data
- Identify groups in the data

Greater similarity between data points results in better clustering.

# Unsupervised Learning Models: K-Means Clustering

K-means finds the best centroids by alternatively assigning random centroids to a dataset and selecting mean data points from the resulting clusters to form new centroids. It continues this process iteratively until the model is optimized.



Centroids (mean)

Assign

Find the number of clusters and assign mean

Optimize

Iterate and optimize the mean for each cluster for its respective data points

# Unsupervised Learning Models: K-Means Clustering

K-means finds the best centroids by alternatively assigning random centroids to a dataset and selecting mean data points from the resulting clusters to form new centroids. It continues this process iteratively until the model is optimized.



Assign data points to the centroids

Choose a mean from each cluster as a centroid

Reassign data points to new centroids

Iterate the process till the model is optimized

# Unsupervised Learning Models: K-Means Clustering

Let us see how the k-means algorithm works in Scikit-Learn.

Number of clusters to form and number of centroids to generate

Number of times the K-means algorithm will be run with different centroid seeds

Pre-compute for faster operation

Class

Selects initial cluster centers

sklearn.cluster.KMeans(*n_clusters=8*, *init='k-means++'*, *n_init=10*, *max_iter=300*, *tol=0.0001*, *precompute_distances='auto'*, *verbose=0*, *random_state=None*, *copy_x=True*, *n_jobs=1*)

Initialize the centers

If true, does not modify data while pre-computing

Number of jobs in parallel computation

Maximum number of iterations of the K-means algorithm for a single run

**Problem Statement:** Demonstrate how to use k-means clustering to classify data points.

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

ASSISTED PRACTICE

**Unsupervised Learning Models: Dimensionality Reduction**

# Unsupervised Learning Models: Dimensionality Reduction

It reduces a high-dimensional dataset into a dataset with fewer dimensions. This makes it easier and faster for the algorithm to analyze the data.

# Unsupervised Learning Models: Dimensionality Reduction

These are some techniques used for dimensionality reduction:



Large dataset
(a few thousand columns and rows)

Drop data columns with missing values

Drop data columns with low variance

Drop data columns with high correlations

Apply statistical functions - PCA

Unsupervised Learning Models: Principal Component Analysis

# Unsupervised Learning Models: Principal Component Analysis (PCA)

It is a linear dimensionality reduction method which uses singular value decomposition of the data and keeps only the most significant singular vectors to project the data to a lower dimensional space.

- It is primarily used to compress or reduce the data.

- PCA tries to capture the variance, which helps it pick up interesting features.

- PCA is used to reduce dimensionality in the dataset and to build our feature vector.

- Here, the principal axes in the feature space represents the direction of maximum variance in the data.

Minor axis

Principal axis

This method is used to capture variance.

# Unsupervised Learning Models: Principal Component Analysis

Let us look at how the PCA algorithm works in Scikit-Learn.

Number of components to keep

Class

sklearn.decomposition.PCA(*n_components=None, copy=True, whiten=False*)

Overwrites the transform data after fitting them into the model

Removes data with lower variance

# Principal Component Analysis (PCA)

**Problem Statement:** Demonstrate how to use the PCA model to reduce the dimensions of a dataset.

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

ASSISTED PRACTICE

# Pipeline

It simplifies the process where more than one model is required or used.

All models in the pipeline must be transformers. The last model can either be a transformer or a classifier, regressor, or other such objects.

Once all the data is fit into the models or estimators, the predict method can be called.

Estimators are known as 'model instance'.

# Pipelines

**Problem Statement:** Demonstrate how to build a pipeline.

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

ASSISTED PRACTICE

# Model Persistence

You can save your model for future use. This avoids the need to retrain the model.

- This can be saved using the Pickle method.

- It can also be replaced with the joblib of Sci-kit team.

- Both joblib.dump and joblib.load  can be used.

- These would be efficient for Big Data.

# Model Persistence

**Problem Statement:** Demonstrate how to persist a model for future use.

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

# Model Evaluation: Metric Functions

You can use the "Metrics" function to evaluate the accuracy of your model's predictions.

**Classification** ➡ `metrics.` **accuracy_score**
`metrics.` **average_precision_score**

**Clustering** ➡ `metrics.` **adjusted_rand_score**

**Regression** ➡ `metrics.` **mean_absolute_error**
`metrics.` **mean_squared_error**
`metrics.` **median_absolute_error**

# Evaluation of Dataset

**Problem Statement:**

The given dataset contains ad budgets for different media channels and the corresponding ad sales of XYZ firm. Evaluate the dataset to:
- Find the features or media channels used by the firm
- Find the sales figures for each channel
- Create a model to predict the sales outcome
- Split as training and testing datasets for the model
- Calculate the Mean Square Error (MSE)

**Instructions to perform the assignment:**
•Download the FAA dataset from the "Resource" tab. Upload the dataset to your Jupyter notebook to view and evaluate it.

**Common instructions:**
•If you are new to Python, download the "Anaconda Installation Instructions" document from the "Resources" tab to view the steps for installing Anaconda and the Jupyter notebook.
•Download the "Assignment 01" notebook and upload it on the Jupyter notebook to access it.
•Follow the cues provided to complete the assignment.

**Problem Statement:**

The given dataset lists the glucose level readings of several pregnant women taken either during a survey examination or routine medical care. It specifies if the two hours post-load plasma glucose was at least 200 mg/dl. Analyze the dataset to:
1. Find the features of the dataset
2. Find the response label of the dataset
3. Create a model to predict the diabetes outcome
4. Use training and testing datasets to train the model
5. Check the accuracy of the model

# Listing the Glucose Level Readings

**Instructions on performing the assignment:**
•Download the "pima-Indians-diabetes.DATA" and "pima-Indians-diabetes.NAMES" files from the "Resources" tab. Load the .DATA file to the Jupyter notebook to work on it.
•Open the .NAMES file with a notepad application to view its text. Use this file to view the features of the dataset and add them manually in your code.

**Common instructions:**
•If you are new to Python, download the "Anaconda Installation Instructions" document from the "Resources" tab to view the steps for installing Anaconda and the Jupyter notebook.
•Download the "Assignment 02" notebook and upload it on the Jupyter notebook to access it.
•Follow the provided cues to complete the assignment.

UNASSISTED PRACTICE

# Key Takeaways

You are now able to:

- Define machine learning

- Explain the machine learning approach

- List relevant terminologies that help you understand a dataset

- Discuss features of supervised and unsupervised learning models

- Explain algorithms, such as regression, classification, clustering, and dimensionality reduction

Knowledge Check

**In machine learning, which one of the following is an observation?**

a. Features

b. Attributes

c. Records

d. Labels

**In machine learning, which one of the following is an observation?**

a. Features

b. Attributes

c. Records

d. Labels

The correct answer is **C**

**An observation is a set of examples, records, or samples.**

**Knowledge Check**

**2**

**If data is continuous and has labels (response), then it fits which of the following problem types?**

a. Supervised learning: Classification

b. Unsupervised learning: Clustering

c. Unsupervised learning: Dimensionality reduction

d. Supervised learning: Regression

**If data is continuous and has labels (response), then it fits which of the following problem types?**

a.   Supervised learning: Classification

b.   Unsupervised learning: Clustering

c.   Unsupervised learning: Dimensionality reduction

d.   Supervised learning: Regression

The correct answer is   **d**

**The regression algorithm belonging to the supervised learning model is best suited to analyze continuous data.**

**Identify the goal of unsupervised learning. Select all that apply.**

a. To predict the outcome

b. To understand the structure of the data

c. To generalize the dataset

d. To represent the data

**Identify the goal of unsupervised learning. Select all that apply.**

a. To predict the outcome

b. To understand the structure of the data

c. To generalize the dataset

d. To represent the data

The correct answer is **b, d**

**The goal of unsupervised learning is to understand the structure of the data and represent it. There is no right or certain answer in unsupervised learning.**

**Knowledge Check**

**4**

**The estimator instance in Scikit-learn is a _____.**

a.  To predict the outcome

b.  To understand the structure of the data

c.  To generalize the dataset

d.  To represent the data

**Knowledge Check**

**4**

**The estimator instance in Scikit-learn is a _____.**

a. Model

b. Feature

c. Dataset

d. Response

The correct answer is  **a**

**The estimator instance or object is a model.**

**What is the best way to train a model?**

a. Use the entire dataset as a training and testing set both

b. Split the known dataset into separate training and testing sets

c. Ask the source to provide continuous data

d. Ask the source to provide categorical data

## Knowledge Check 5

**What is the best way to train a model?**

a. Use the entire dataset as a training and testing set both

b. Split the known dataset into separate training and testing sets

c. Ask the source to provide continuous data

d. Ask the source to provide categorical data

The correct answer is **b**

**The best way to train a model is to split the known dataset into training and testing sets. The testing set varies from 20% to 40%.**

**Which of the following is true with a greater value of SSR or SSE? Select all that apply.**

a.  The prediction will be more accurate, making it the best fit model.

b.  The prediction will start becoming less accurate.

c.  The outcome remains unaffected.

d.  The model will not be the best fit for the attributes.

**Knowledge Check 6**

**Which of the following is true with a greater value of SSR or SSE? Select all that apply.**

a.  The prediction will be more accurate, making it the best fit model.

b.  The prediction will start becoming less accurate.

c.  The outcome remains unaffected.

d.  The model will not be the best fit for the attributes.

The correct answer is  **b, d**

**With higher SSR or SSE, the prediction will be less accurate and the model will not be the best fit for the attributes.**

**Class sklearn.linear_model.LogisticRegression, random_state _____.**

a. Indicates the seed of the pseudo random number generator used to shuffle data

b. Defines the features state

c. Represents the number of random iterations

d. Specifies a random constant to be added to the decision function

**Class sklearn.linear_model.LogisticRegression, random_state _____.**

a. Indicates the seed of the pseudo random number generator used to shuffle data

b. Defines the features state

c. Represents the number of random iterations

d. Specifies a random constant to be added to the decision function

The correct answer is **a**

**The class "sklearn.linear_model.LogisticRegression, random_state" indicates the seed of the pseudo random number generator used to shuffle data.**

**What are the requirements of the K-means algorithm? Select all that apply.**

a. Number of clusters should be specified

b. More than one iteration should meet requisite criteria

c. Centroids should minimize inertia

d. Features should be labeled

**Knowledge Check**

**8**

**What are the requirements of the K-means algorithm? Select all that apply.**

a. Number of clusters should be specified

b. More than one iteration should meet requisite criteria

c. Centroids should minimize inertia

d. Features should be labeled

The correct answer is **a, b, c**

**The K-means algorithm requires the number of clusters to be specified and the centroids to minimize inertia. It requires several iterations to fine tune itself and meet the required criteria to become the best fit model.**

**What are the requirements of the K-means algorithm? Select all that apply.**

a.     Fits the model with X and applies the dimensionality reduction on X

b.     Transforms the data back to its original space

c.     Applies the dimensionality reduction on X

d.     Computes data co-variance with the generative model

simplilearn

**In Class sklearn.decomposition.PCA, the transform(X) method , where X is multi-dimensional _____.**

a. Fits the model with X and applies the dimensionality reduction on X

b. Transforms the data back to its original space

c. Applies the dimensionality reduction on X

d. Computes data co-variance with the generative model

The correct answer is   **C**

**In Class "sklearn.decomposition.PCA," the transform(X) method applies the dimensionality reduction on X.**