

# Unit II

**Applications of machine learning in Data science, role of ML in DS, Python tools like sklearn, modelling process for feature engineering, model selection, validation and prediction, types of ML, semi-supervised learning**

**Handling large data: problems and general techniques for handling large data, programming tips for dealing large data, case studies on DS projects for predicting malicious URLs, for building recommender systems**

## 3.1. What is machine learning?

“Machine learning is a field of study that gives computers the ability to learn without being explicitly programmed.”

Or

Machine learning (ML) is a branch of [Artificial Intelligence \(AI\)](#) and computer science that focuses on the using data and algorithms to enable AI to imitate the way that humans learn, gradually improving its accuracy.

Or

Machine learning (ML) is a discipline of artificial intelligence (AI) that provides machines with the ability to automatically learn from data and past experiences while identifying patterns to make predictions with minimal human intervention.

### Data Science vs. Machine Learning

DATA SCIENCE	MACHINE LEARNING
Data Science is all about processes and systems to extract data from structured and semi-structured data.	Machine learning (ML) is a field of artificial intelligence (AI) that allows the software to learn from data to identify patterns and make predictions automatically with minimal human intervention.
Data Science deals with data.	Machine learning utilizes past experiences to learn about the data.
Data gathering, manipulation, Data cleaning, etc., are data science operations.	ML is of three types: Supervised learning, Unsupervised learning, and Reinforcement learning.
Example: Netflix using data Science is an example of this technology.  With the help of Data Science, Netflix can provide users with personalized recommendations on movies and shows. It can also predict the original content's popularity with trailers and thumbnail images.	Example: Facebook using ML is an example of this technology.  With the help of Machine learning, Facebook can produce the estimated action rate and the ad quality score, which is used for the total equation. ML features such as facial recognition, targeted advertising, language

	translation, and news ed are also used in many real-case scenarios.
--	---

### 3.1.1. Applications for machine learning in data science

*Regression* and *classification* are of primary importance to a data scientist. To achieve these goals, one of the main tools a data scientist uses is machine learning. The uses for regression and automatic classification are wide ranging, such as the following:

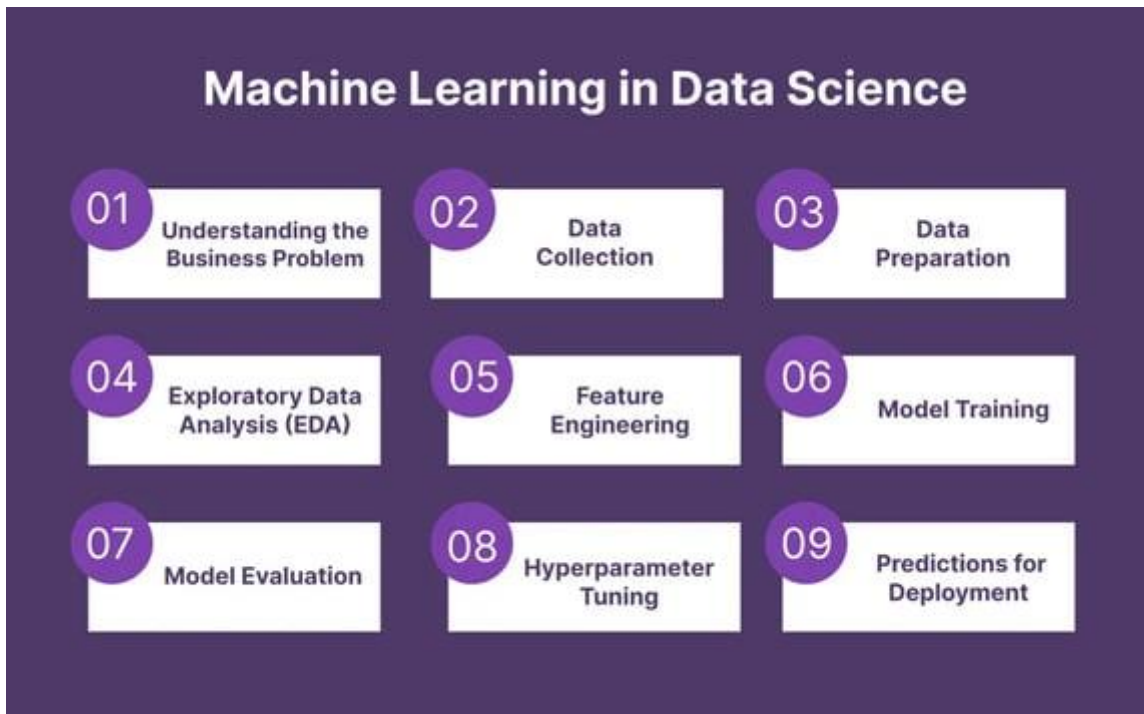
- Finding oil fields, gold mines, or archeological sites based on existing sites (classification and regression)
- Finding place names or persons in text (classification)
- Identifying people based on pictures or voice recordings (classification)
- Recognizing birds based on their whistle (classification)
- Identifying profitable customers (regression and classification)
- Proactively identifying car parts that are likely to fail (regression)
- Identifying tumors and diseases (classification)
- Predicting the amount of money a person will spend on product X (regression)
- Predicting the number of eruptions of a volcano in a period (regression)
- Predicting your company's yearly revenue (regression)
- Predicting which team will win the Champions League in soccer (classification)

Occasionally data scientists build a *model* (an abstraction of reality) that provides insight to the underlying processes of a phenomenon. When the goal of a model isn't prediction but interpretation, it's called *root cause analysis*. Here are a few examples:

- Understanding and optimizing a business process, such as determining which products add value to a product line
- Discovering what causes diabetes
- Determining the causes of traffic jams

This list of machine learning applications can only be seen as an appetizer because it's ubiquitous within data science. Regression and classification are two important techniques, but the repertoire and the applications don't end, with clustering as one other example of a valuable technique. Machine learning techniques can be used throughout the data science process, as we'll discuss in the next section.

## The Role of Machine Learning in Data Science



The role of machine learning in Data Science occurs in 9 steps:

### 1. Understanding the Business Problem

To build a successful business model, it's very important to understand the business problem that the client is facing. Suppose the client wants to predict whether the patient has cancer or not. In such scenario, domain experts understand the underlying problems that are present in the system.

### 2. Data Collection

After understanding the problem statement, you have to collect relevant data. As per the business problem, machine learning helps collect and analyze structured, unstructured, and semi-structured data from any database across systems.

### 3. Data Preparation

The first step of data preparation is data cleaning. It is an essential step for preparing the data. In data preparation, you eliminate duplicates and null values, inconsistent data types, invalid entries, missing data, and improper formatting.

### 4. Exploratory Data Analysis (EDA)

Exploratory Data Analysis lets you uncover valuable insights that will be useful in the next phase of the Data Science lifecycle. EDA is important because, through EDA, you can find outliers, anomalies, and trends in the dataset. These insights can be helpful in identifying the optimal number of features to be used for model building.

### 5. Feature Engineering

Feature engineering is one of the important steps in a Data Science Project. It helps in creating new features, transforming and scaling the features. In this domain, expertise plays a key role in generating new insights from the data exploration step.

### 6. Model Training

In Model training, we fit the training data; this is where “learning” starts. We train the model on training data and test the performance on testing data i.e., unseen data.

## 7. Model Evaluation

Once Model Training is done, it's time to evaluate its performance. So, evaluating your Model on a new dataset will give you an idea of how your Model is going to perform in future data.

## 8. Hyperparameter Tuning

After the Model is trained and evaluated, the performance of the Model can be again improved by tuning its parameter. Hyperparameter tuning of the model is important to improve the overall performance of the model.

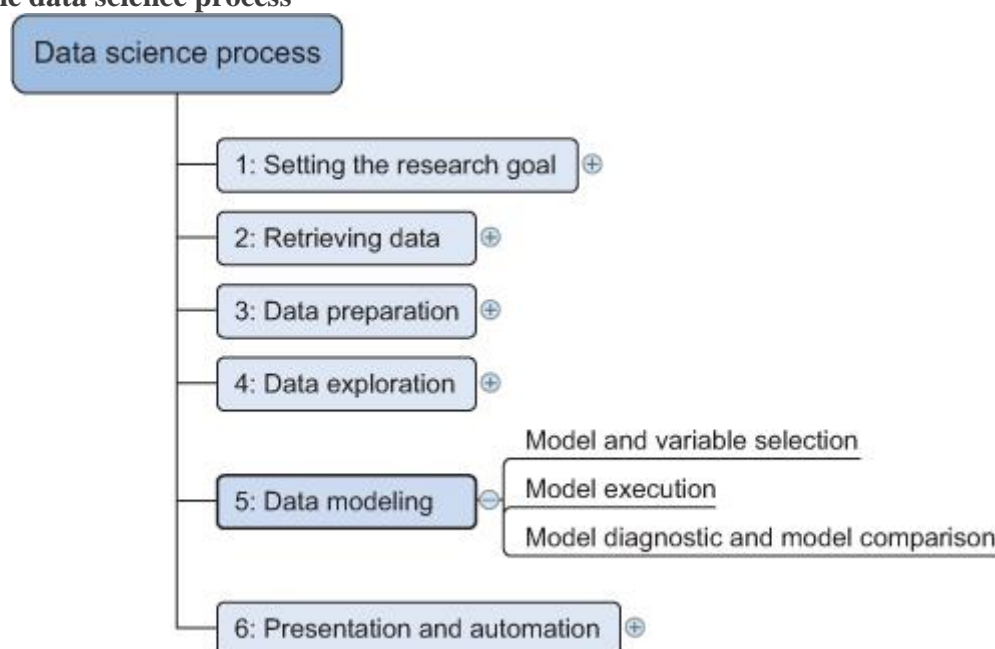
## 9. Making Predictions and Ready to be Deployed

This is the final stage of machine learning. Here, the machine answers each of your questions by its learning. After making accurate predictions, the Data Model is deployed into production.

### 3.1.2. Where machine learning is used in the data science process

Although machine learning is mainly linked to the data-modeling step of the data science process, it can be used at almost every step.

**Figure 3.1. The data science process**



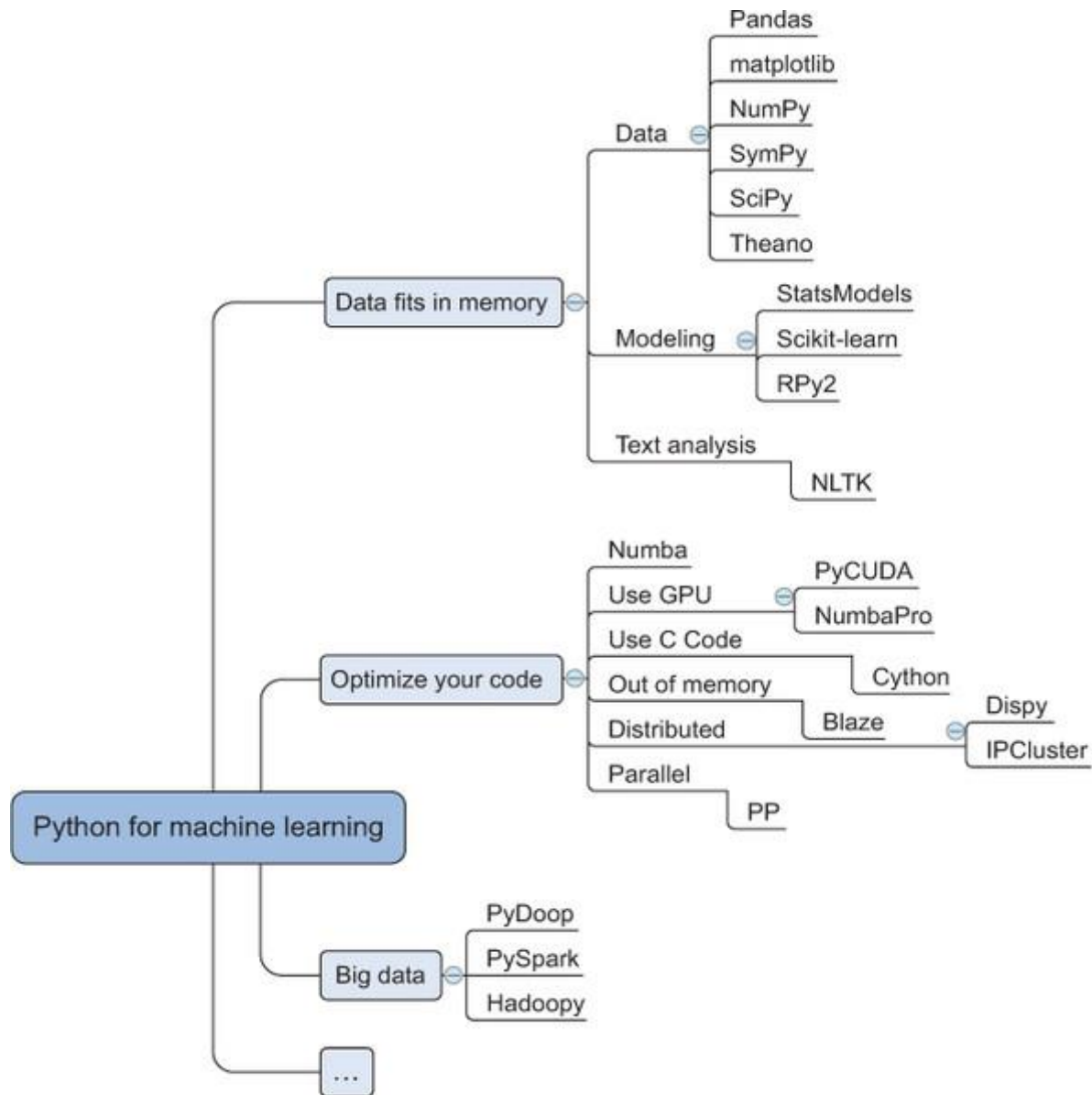
The data modeling phase can't start until you have qualitative raw data you can understand. But prior to that, the *data preparation* phase can benefit from the use of machine learning. An example would be cleansing a list of text strings; machine learning can group similar strings together so it becomes easier to correct spelling errors. Machine learning is also useful when *exploring data*. Algorithms can root out underlying patterns in the data where they'd be difficult to find with only charts.

Given that machine learning is useful throughout the data science process, it shouldn't come as a surprise that a considerable number of Python libraries were developed to make your life a bit easier.

### 3.1.3. Python tools used in machine learning

Python has an overwhelming number of packages that can be used in a machine learning setting. The Python machine learning ecosystem can be divided into three main types of packages, as shown in [figure 3.2](#).

Figure 3.2. Overview of Python packages used during the machine-learning phase



The first type of package shown in [figure 3.2](#) is mainly used in simple tasks and when data fits into memory. The second type is used to optimize your code when you've finished prototyping and run into speed or memory issues. The third type is specific to using Python with big data technologies.

#### Packages for working with data in memory

When prototyping, the following packages can get you started by providing advanced functionalities with a few lines of code:

- **SciPy** is a library that integrates fundamental packages often used in scientific computing such as NumPy, matplotlib, Pandas, and SymPy.
- **NumPy** gives you access to powerful array functions and linear algebra functions.
- **Matplotlib** is a popular 2D plotting package with some 3D functionality.

- **Pandas** is a high-performance, but easy-to-use, data-wrangling package. It introduces dataframes to Python, a type of in-memory data table. It's a concept that should sound familiar to regular users of R.
- **SymPy** is a package used for symbolic mathematics and computer algebra.
- **StatsModels** is a package for statistical methods and algorithms.
- **Scikit-learn** is a library filled with machine learning algorithms.
- **RPy2** allows you to call R functions from within Python. R is a popular open source statistics program.
- **NLTK** (Natural Language Toolkit) is a Python toolkit with a focus on text analytics.

These libraries are good to get started with, but once you make the decision to run a certain Python program at frequent intervals, performance comes into play.

### Optimizing operations

Once your application moves into production, the libraries listed here can help you deliver the speed you need. Sometimes this involves connecting to big data infrastructures such as Hadoop and Spark.

- **Numba and NumbaPro** —These use just-in-time compilation to speed up applications written directly in Python and a few annotations. NumbaPro also allows you to use the power of your graphics processor unit (GPU).
- **PyCUDA** —This allows you to write code that will be executed on the GPU instead of your CPU and is therefore ideal for calculation-heavy applications. It works best with problems that lend themselves to being parallelized and need little input compared to the number of required computing cycles. An example is studying the robustness of your predictions by calculating thousands of different outcomes based on a single start state.
- **Cython, or C for Python** —This brings the C programming language to Python. C is a lower-level language, so the code is closer to what the computer eventually uses (bytecode). The closer code is to bits and bytes, the faster it executes. A computer is also faster when it knows the type of a variable (called *static typing*). Python wasn't designed to do this, and Cython helps you to overcome this shortfall.
- **Blaze** —Blaze gives you data structures that can be bigger than your computer's main memory, enabling you to work with large data sets.
- **Dispy and IPCluster** —These packages allow you to write code that can be distributed over a cluster of computers.
- **PP** —Python is executed as a single process by default. With the help of PP you can parallelize computations on a single machine or over clusters.
- **Pydoop and Hadoopy** —These connect Python to Hadoop, a common big data framework.
- **PySpark** —This connects Python and Spark, an in-memory big data framework.

Now that you've seen an overview of the available libraries, let's look at the modeling process itself.

## 3.2. The modeling process

The modeling phase consists of four steps:

1. **Feature engineering and model selection**
2. **Training the model**
3. **Model validation and selection**
4. **Applying the trained model to unseen data**

Before you find a good model, you'll probably iterate among the first three steps.

The last step isn't always present because sometimes the goal isn't prediction but explanation (root cause analysis). For instance, you might want to find out the causes of species' extinctions but not necessarily predict which one is next in line to leave our planet.

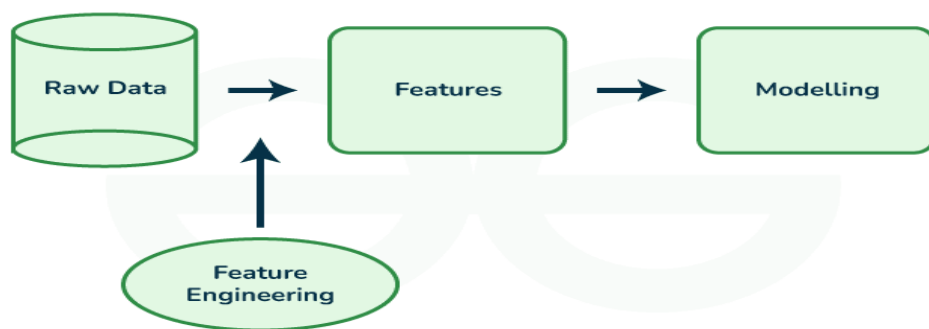
It's possible to *chain* or *combine* multiple techniques. When you chain multiple models, the output of the first model becomes an input for the second model. When you combine multiple models, you train them independently and combine their results. This last technique is also known as *ensemble learning*.

A model consists of constructs of information called *features* or *predictors* and a *target* or *response variable*. Your model's goal is to predict the target variable, for example, tomorrow's high temperature. The variables that help you do this and are (usually) known to you are the features or predictor variables such as today's temperature, cloud movements, current wind speed, and so on. The best models are those that accurately represent reality, preferably while staying concise and interpretable. To achieve this, feature engineering is the most important and arguably most interesting part of modeling. For example, an important feature in a model that tried to explain the extinction of large land animals in the last 60,000 years in Australia turned out to be the population number and spread of humans.

### 3.2.1. Engineering features and selecting a model

Feature engineering is the process of transforming raw data into features that are suitable for machine learning models. In other words, it is the process of selecting, extracting, and transforming the most relevant features from the available data to build more accurate and efficient machine learning models.

The success of machine learning models heavily depends on the quality of the features used to train them. Feature engineering involves a set of techniques that enable us to create new features by combining or transforming the existing ones. These techniques help to highlight the most important patterns and relationships in the data, which in turn helps the machine learning model to learn from the data more effectively.



### 3.2.2. Training your model

In this phase you present to your model data from which it can learn.

The most common modeling techniques have industry-ready implementations in almost every programming language, including Python. These enable you to train your models by executing a few lines of code. For more state-of-the-art data science techniques, you'll probably end up doing heavy mathematical calculations and implementing them with modern computer science techniques.

Once a model is trained, it's time to test whether it can be extrapolated to reality: model validation.

### 3.2.3. Validating a model

Data science has many modeling techniques, and the question is which one is the right one to use. A good model has two properties: it has good predictive power and it generalizes well to data it hasn't seen. To achieve this you define an error measure (how wrong the model is) and a validation strategy.

Two common *error measures* in machine learning are the *classification error rate* for classification problems and the *mean squared error* for regression problems. The classification error rate is the percentage of observations in the test data set that your model mislabeled; lower is better. The mean squared error measures how big the average error of your prediction is. Squaring the average error has two consequences: you can't

cancel out a wrong prediction in one direction with a faulty prediction in the other direction. For example, overestimating future turnover for next month by 5,000 doesn't cancel out underestimating it by 5,000 for the following month. As a second consequence of squaring, bigger errors get even more weight than they otherwise would. Small errors remain small or can even shrink (if  $<1$ ), whereas big errors are enlarged and will definitely draw your attention.

Many *validation strategies* exist, including the following common ones:

- *Dividing your data into a training set with X% of the observations and keeping the rest as a holdout data set* (a data set that's never used for model creation)—This is the most common technique.
- **K-folds cross validation** —This strategy divides the data set into k parts and uses each part one time as a test data set while using the others as a training data set. This has the advantage that you use all the data available in the data set.
- **Leave-1 out** —This approach is the same as k-folds but with  $k=1$ . You always leave one observation out and train on the rest of the data. This is used only on small data sets, so it's more valuable to people evaluating laboratory experiments than to big data analysts.

#### 3.2.4. Predicting new observations

If you've implemented the first three steps successfully, you now have a performant model that generalizes to unseen data. The process of applying your model to new data is called model scoring. In fact, model scoring is something you implicitly did during validation, only now you don't know the correct outcome. By now you should trust your model enough to use it for real.

Model scoring involves two steps. First, you prepare a data set that has features exactly as defined by your model. This boils down to repeating the data preparation you did in step one of the modeling process but for a new data set. Then you apply the model on this new data set, and this results in a prediction.

Now let's look at the different types of machine learning techniques: a different problem requires a different approach.

### 3.3. Types of machine learning

## Types of Machine Learning

**Machine learning is a subset of AI, which enables the machine to automatically learn from data, improve performance from past experiences, and make predictions.** Machine learning contains a set of algorithms that work on a huge amount of data. Data is fed to these algorithms to train them, and on the basis of training, they build the model & perform a specific task.



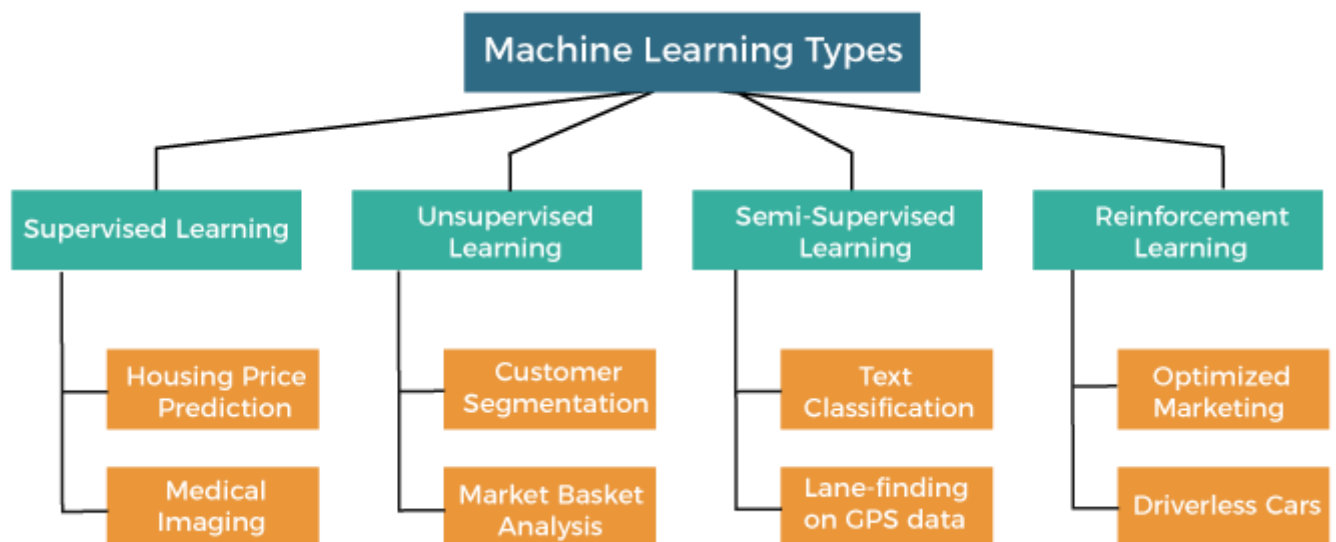
# Types of Machine Learning



These ML algorithms help to solve different business problems like Regression, Classification, Forecasting, Clustering, and Associations, etc.

Based on the methods and way of learning, machine learning is divided into mainly four types, which are:

1. Supervised Machine Learning
2. Unsupervised Machine Learning
3. Semi-Supervised Machine Learning
4. Reinforcement Learning



In this topic, we will provide a detailed description of the types of Machine Learning along with their respective algorithms:

## 1. Supervised Machine Learning

As its name suggests, [Supervised machine learning](#) is based on supervision. It means in the supervised learning technique, we train the machines using the "labelled" dataset, and based on the training, the machine predicts the output. Here, the labelled data specifies that some of the inputs are already mapped to the output. More precisely, we can say; first, we train the machine with the input and corresponding output, and then we ask the machine to predict the output using the test dataset.

Let's understand supervised learning with an example. Suppose we have an input dataset of cats and dog images. So, first, we will provide the training to the machine to understand the images, such as the **shape & size of the tail of cat and dog, Shape of eyes, colour, height (dogs are taller, cats are smaller), etc.** After completion of training, we input the picture of a cat and ask the machine to identify the object and predict the output. Now, the machine is well trained, so it will check all the features of the object, such as height, shape, colour, eyes, ears, tail, etc., and find that it's a cat. So, it will put it in the Cat category. This is the process of how the machine identifies the objects in Supervised Learning.

**The main goal of the supervised learning technique is to map the input variable(x) with the output variable(y).** Some real-world applications of supervised learning are **Risk Assessment, Fraud Detection, Spam filtering**, etc.

## Categories of Supervised Machine Learning

Supervised machine learning can be classified into two types of problems, which are given below

- **Classification**
- **Regression**

### a) Classification

Classification algorithms are used to solve the classification problems in which the output variable is categorical, such as "**Yes" or No, Male or Female, Red or Blue, etc.** The classification algorithms predict the categories present in the dataset. Some real-world examples of classification algorithms are **Spam Detection, Email filtering, etc.**

Some popular classification algorithms are given below:

- **Random Forest Algorithm**
- **Decision Tree Algorithm**
- **Logistic Regression Algorithm**
- **Support Vector Machine Algorithm**

### b) Regression

Regression algorithms are used to solve regression problems in which there is a linear relationship between input and output variables. These are used to predict continuous output variables, such as market trends, weather prediction, etc.

Some popular Regression algorithms are given below:

- **Simple Linear Regression Algorithm**
- **Multivariate Regression Algorithm**
- **Decision Tree Algorithm**
- **Lasso Regression**

## Advantages and Disadvantages of Supervised Learning

### Advantages:

- Since supervised learning work with the labelled dataset so we can have an exact idea about the classes of objects.

- These algorithms are helpful in predicting the output on the basis of prior experience.

#### **Disadvantages:**

- These algorithms are not able to solve complex tasks.
- It may predict the wrong output if the test data is different from the training data.
- It requires lots of computational time to train the algorithm.

## Applications of Supervised Learning

Some common applications of Supervised Learning are given below:

- **Image Segmentation:**  
Supervised Learning algorithms are used in image segmentation. In this process, image classification is performed on different image data with pre-defined labels.
- **Medical Diagnosis:**  
Supervised algorithms are also used in the medical field for diagnosis purposes. It is done by using medical images and past labelled data with labels for disease conditions. With such a process, the machine can identify a disease for the new patients.
- **Fraud Detection** - Supervised Learning classification algorithms are used for identifying fraud transactions, fraud customers, etc. It is done by using historic data to identify the patterns that can lead to possible fraud.
- **Spam detection** - In spam detection & filtering, classification algorithms are used. These algorithms classify an email as spam or not spam. The spam emails are sent to the spam folder.
- **Speech Recognition** - Supervised learning algorithms are also used in speech recognition. The algorithm is trained with voice data, and various identifications can be done using the same, such as voice-activated passwords, voice commands, etc.

## 2. Unsupervised Machine Learning

[Unsupervised learning](#) is different from the Supervised learning technique; as its name suggests, there is no need for supervision. It means, in unsupervised machine learning, the machine is trained using the unlabeled dataset, and the machine predicts the output without any supervision.

In unsupervised learning, the models are trained with the data that is neither classified nor labelled, and the model acts on that data without any supervision.

**The main aim of the unsupervised learning algorithm is to group or categories the unsorted dataset according to the similarities, patterns, and differences.** Machines are instructed to find the hidden patterns from the input dataset.

Let's take an example to understand it more preciously; suppose there is a basket of fruit images, and we input it into the machine learning model. The images are totally unknown to the model, and the task of the machine is to find the patterns and categories of the objects.

So, now the machine will discover its patterns and differences, such as colour difference, shape difference, and predict the output when it is tested with the test dataset.

## Categories of Unsupervised Machine Learning

Unsupervised Learning can be further classified into two types, which are given below:

- **Clustering**

- **Association**

## 1) Clustering

The clustering technique is used when we want to find the inherent groups from the data. It is a way to group the objects into a cluster such that the objects with the most similarities remain in one group and have fewer or no similarities with the objects of other groups. An example of the clustering algorithm is grouping the customers by their purchasing behaviour.

Some of the popular clustering algorithms are given below:

- **K-Means Clustering algorithm**
- **Mean-shift algorithm**
- **DBSCAN Algorithm**
- **Principal Component Analysis**
- **Independent Component Analysis**

## 2) Association

Association rule learning is an unsupervised learning technique, which finds interesting relations among variables within a large dataset. The main aim of this learning algorithm is to find the dependency of one data item on another data item and map those variables accordingly so that it can generate maximum profit. This algorithm is mainly applied in **Market Basket analysis, Web usage mining, continuous production**, etc.

Some popular algorithms of Association rule learning are **Apriori Algorithm, Eclat, FP-growth algorithm**.

## Advantages and Disadvantages of Unsupervised Learning Algorithm

### Advantages:

- These algorithms can be used for complicated tasks compared to the supervised ones because these algorithms work on the unlabeled dataset.
- Unsupervised algorithms are preferable for various tasks as getting the unlabeled dataset is easier as compared to the labelled dataset.

### Disadvantages:

- The output of an unsupervised algorithm can be less accurate as the dataset is not labelled, and algorithms are not trained with the exact output in prior.
- Working with Unsupervised learning is more difficult as it works with the unlabelled dataset that does not map with the output.

## Applications of Unsupervised Learning

- **Network Analysis:** Unsupervised learning is used for identifying plagiarism and copyright in document network analysis of text data for scholarly articles.

- **Recommendation Systems:** Recommendation systems widely use unsupervised learning techniques for building recommendation applications for different web applications and e-commerce websites.
- **Anomaly Detection:** Anomaly detection is a popular application of unsupervised learning, which can identify unusual data points within the dataset. It is used to discover fraudulent transactions.
- **Singular Value Decomposition:** Singular Value Decomposition or SVD is used to extract particular information from the database. For example, extracting information of each user located at a particular location.

### 3. Semi-Supervised Learning

**Semi-Supervised learning is a type of Machine Learning algorithm that lies between Supervised and Unsupervised machine learning.** It represents the intermediate ground between Supervised (With Labelled training data) and Unsupervised learning (with no labelled training data) algorithms and uses the combination of labelled and unlabeled datasets during the training period.

Although Semi-supervised learning is the middle ground between supervised and unsupervised learning and operates on the data that consists of a few labels, it mostly consists of unlabeled data. As labels are costly, but for corporate purposes, they may have few labels. It is completely different from supervised and unsupervised learning as they are based on the presence & absence of labels.

**To overcome the drawbacks of supervised learning and unsupervised learning algorithms, the concept of Semi-supervised learning is introduced.** The main aim of [semi-supervised learning](#) is to effectively use all the available data, rather than only labelled data like in supervised learning. Initially, similar data is clustered along with an unsupervised learning algorithm, and further, it helps to label the unlabeled data into labelled data. It is because labelled data is a comparatively more expensive acquisition than unlabeled data.

We can imagine these algorithms with an example. Supervised learning is where a student is under the supervision of an instructor at home and college. Further, if that student is self-analysing the same concept without any help from the instructor, it comes under unsupervised learning. Under semi-supervised learning, the student has to revise himself after analyzing the same concept under the guidance of an instructor at college.

### Advantages and disadvantages of Semi-supervised Learning

#### Advantages:

- It is simple and easy to understand the algorithm.
- It is highly efficient.
- It is used to solve drawbacks of Supervised and Unsupervised Learning algorithms.

#### Disadvantages:

- Iterations results may not be stable.
- We cannot apply these algorithms to network-level data.
- Accuracy is low.

### 4. Reinforcement Learning

**Reinforcement learning works on a feedback-based process, in which an AI agent (A software component) automatically explore its surrounding by hitting & trail, taking action, learning from experiences, and improving its performance.** Agent gets rewarded for each good action and get punished for each bad action; hence the goal of reinforcement learning agent is to maximize the rewards.

In reinforcement learning, there is no labelled data like supervised learning, and agents learn from their experiences only.

The [reinforcement learning](#) process is similar to a human being; for example, a child learns various things by experiences in his day-to-day life. An example of reinforcement learning is to play a game, where the Game is the environment, moves of an agent at each step define states, and the goal of the agent is to get a high score. Agent receives feedback in terms of punishment and rewards.

Due to its way of working, reinforcement learning is employed in different fields such as **Game theory, Operation Research, Information theory, multi-agent systems.**

A reinforcement learning problem can be formalized using **Markov Decision Process(MDP)**. In MDP, the agent constantly interacts with the environment and performs actions; at each action, the environment responds and generates a new state.

## Categories of Reinforcement Learning

Reinforcement learning is categorized mainly into two types of methods/algorithms:

- **Positive Reinforcement Learning:** Positive reinforcement learning specifies increasing the tendency that the required behaviour would occur again by adding something. It enhances the strength of the behaviour of the agent and positively impacts it.
- **Negative Reinforcement Learning:** Negative reinforcement learning works exactly opposite to the positive RL. It increases the tendency that the specific behaviour would occur again by avoiding the negative condition.

## Real-world Use cases of Reinforcement Learning

- **Video Games:**  
RL algorithms are much popular in gaming applications. It is used to gain super-human performance. Some popular games that use RL algorithms are **AlphaGO** and **AlphaGO Zero**.
- **Resource Management:**  
The "Resource Management with Deep Reinforcement Learning" paper showed that how to use RL in computer to automatically learn and schedule resources to wait for different jobs in order to minimize average job slowdown.
- **Robotics:**  
RL is widely being used in Robotics applications. Robots are used in the industrial and manufacturing area, and these robots are made more powerful with reinforcement learning. There are different industries that have their vision of building intelligent robots using AI and Machine learning technology.
- **Text Mining**  
Text-mining, one of the great applications of NLP, is now being implemented with the help of Reinforcement Learning by Salesforce company.

## Advantages and Disadvantages of Reinforcement Learning

### Advantages

- It helps in solving complex real-world problems which are difficult to be solved by general techniques.
- The learning model of RL is similar to the learning of human beings; hence most accurate results can be found.
- Helps in achieving long term results.

#### **Disadvantage**

- RL algorithms are not preferred for simple problems.
- RL algorithms require huge data and computations.
- Too much reinforcement learning can lead to an overload of states which can weaken the results.

The curse of dimensionality limits reinforcement learning for real physical systems.

**Handling large data on a single computer:** The problems you face when handling large data

**General techniques for handling large volumes of data:** Choosing the right algorithm, Choosing the right data structure, Selecting the right tools

**General programming tips for dealing with large data sets:** Don't reinvent the wheel, Get the most out of your hardware, Reduce your computing needs.

## **Handling large data:**

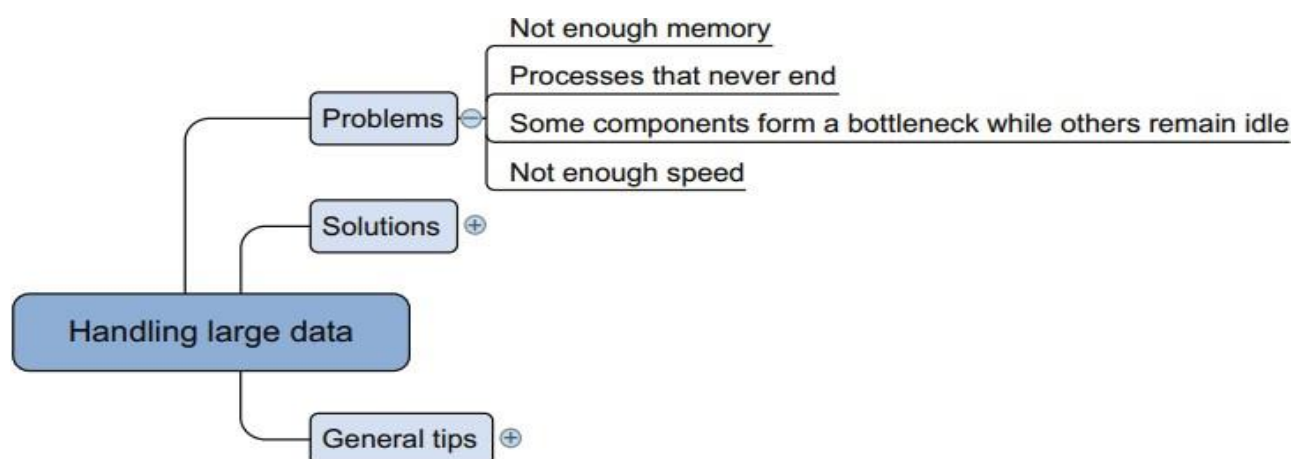
### **Handling large data on a single computer: -**

What if you had so much data that it seems to outgrow you and your techniques no longer seem to suffice? What do you do, surrender or adapt? Luckily you chose to adapt, because you're still reading.

### **The problems you face when handling large data:**

-

A large volume of data poses new challenges, such as overloaded memory and algorithms that never stop running. It forces you to adapt and expand your repertoire of techniques. But even when you can perform your analysis, you should take care of issues such as I/O (input/output) and CPU starvation, because these can cause speed issues. Figure 4.1 shows a mind map that will gradually unfold as we go through the steps: problems, solutions, and tips.



**Figure 4.1 Overview of problems encountered when working with more data than can fit in memory**

A computer only has a limited amount of RAM. When you try to squeeze more data into this memory than actually fits, the OS will start swapping out memory blocks to disks, which is far less efficient than having it all in memory. But only a few algorithms are designed to handle large data sets; most of them load the whole data set into memory at once, which causes the out-of-memory error. Other algorithms need to hold multiple copies of the data in memory or store intermediate results. All of these aggravate the problem.

Even when you cure the memory issues, you may need to deal with another limited resource: **time**. Although a computer may think you live for millions of years, in reality you won't. Certain algorithms don't take time into account; they'll keep running forever. Other algorithms can't end in a reasonable amount of time when they need to process only a few megabytes of data.

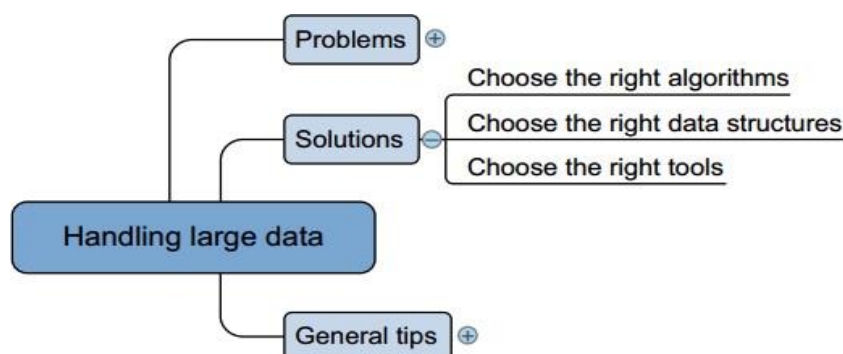
A third thing you'll observe when dealing with large data sets is that components of your computer can start to form a bottleneck while leaving other systems idle. Although this isn't as severe as a never-ending algorithm or out-of-memory errors, it still incurs a serious cost. Think of the cost savings in terms of person days and computing infrastructure for CPU starvation. Certain programs don't feed data fast enough to the processor because they have to read data from the hard drive, which is one of the slowest components on a computer. This has been addressed with the introduction of solid state drives (SSD), but SSDs are still much more expensive than the slower and more widespread hard disk drive (HDD) technology.



# General techniques for handling large volumes of data

Never-ending algorithms, out-of-memory errors, and speed issues are the most common challenges you face when working with large data.

The solutions can be divided into three categories: using the correct algorithms, choosing the right data structure, and using the right tools (figure 4.2).



**Figure 4.2 Overview of solutions for handling large data sets**

No clear one-to-one mapping exists between the problems and solutions because many solutions address both lack of memory and computational performance. For instance, data set compression will help you solve memory issues because the data set becomes smaller. But this also affects computation speed with a shift from the slow hard disk to the fast CPU. Contrary to RAM (random access memory), the hard disc will store everything even after the power goes down, but writing to disc costs more time than changing information in the fleeting RAM. When constantly changing the information, RAM is thus preferable over the (more durable) hard disc. With an unpacked data set, numerous read and write operations (I/O) are occurring, but the CPU remains largely idle, whereas with the compressed data set the CPU gets its fair share of the workload.

## 1. Choosing the right algorithm: -

Choosing the right algorithm can solve more problems than adding more or better hardware. An algorithm that's well suited for handling large data doesn't need to load the entire data set into memory to make predictions. Ideally, the algorithm also supports parallelized calculations. In this section we'll dig into three types of algorithms that can do that: **online algorithms**, **block algorithms**, and **MapReduce algorithms**, as shown in figure 4.3.

---

## ONLINE LEARNING ALGORITHMS: -

Several, but not all, machine learning algorithms can be trained using one observation at a time instead of taking all the data into memory. Upon the arrival of a new data point, the model is trained and the observation can be forgotten; its effect is now incorporated into the model's parameters. For example, a model used to predict the weather can use different parameters (like atmospheric pressure or temperature) in different regions. When the data from one region is loaded into the

algorithm, it forgets about this raw data and moves on to the next region. This “use and forget” way of working is the perfect solution for the memory problem as a single observation is unlikely to ever be big enough to fill up all the

memory of a modern day computer. Listing 4.1 shows how to apply this principle to a **perceptron with online learning**. A perceptron is one of the least complex machine learning algorithms used for binary classification (0 or 1); for instance, will the customer buy or not?

#### Listing 4.1 Training a perceptron by observation

The learning rate of an algorithm is the adjustment it makes every time a new observation comes in. If this is high, the model will adjust quickly to new observations but might “overshoot” and never get precise. An oversimplified example: the optimal (and unknown) weight for an x-variable = 0.75. Current estimation is 0.4 with a learning rate of 0.5; the adjustment =  $0.5 \text{ (learning rate)} * 1 \text{ (size of error)} * 1 \text{ (value of x)} = 0.5$ .  $0.4 \text{ (current weight)} + 0.5 \text{ (adjustment)} = 0.9 \text{ (new weight)}$ , instead of 0.75. The adjustment was too big to get the correct result.

```
import numpy as np
class perceptron():
    def __init__(self, X,y, threshold = 0.5,
learning_rate = 0.1, max_epochs = 10):
        self.threshold = threshold
        self.learning_rate = learning_rate
        self.X = X
        self.y = y
        self.max_epochs = max_epochs
```

Sets up  
perceptron class.

The `__init__` method of any Python class is always run when creating an instance of the class. Several default values are set here.

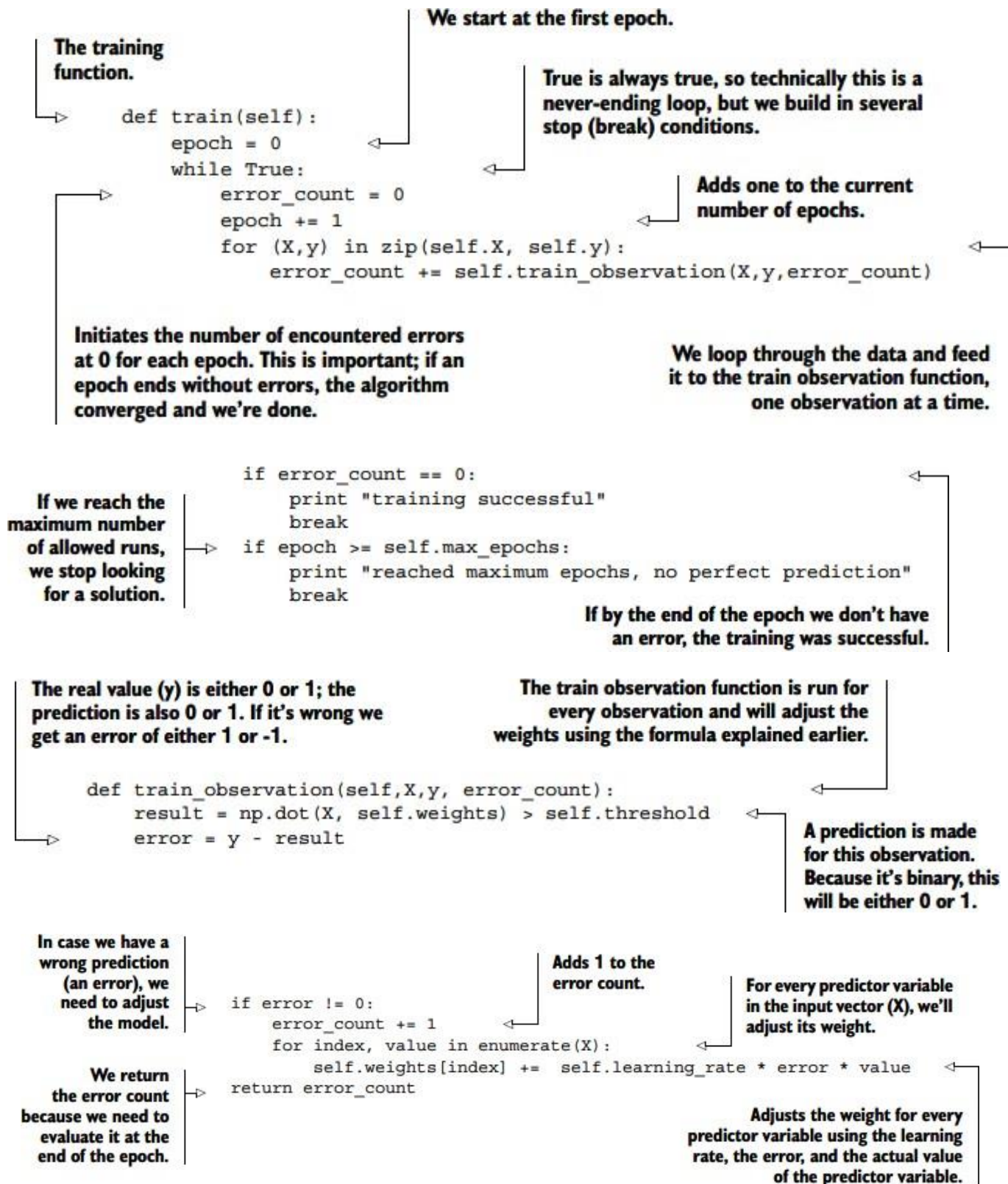
The threshold is an arbitrary cutoff between 0 and 1 to decide whether the prediction becomes a 0 or a 1. Often it's 0.5, right in the middle, but it depends on the use case.

X and y variables are assigned to the class.

One epoch is one run through all the data. We allow for a maximum of 10 runs until we stop the perceptron.

```
def initialize(self, init_type = 'zeros'):
    if init_type == 'random':
        self.weights = np.random.rand(len(self.X[0])) * 0.05
    if init_type == 'zeros':
        self.weights = np.zeros(len(self.X[0]))
```

Each observation will end up with a weight. The initialize function sets these weights for each incoming observation. We allow for 2 options: all weights start at 0 or they are assigned a small (between 0 and 0.05) random weight.





**The predict class.**

```
def predict(self, X):
    return int(np.dot(X, self.weights) > self.threshold)
```

The values of the predictor values are multiplied by their respective weights (this multiplication is done by np.dot). Then the outcome is compared to the overall threshold (here this is 0.5) to see if a 0 or 1 should be predicted.

**Our y (target) data vector.**

```
X = [(1,0,0), (1,1,0), (1,1,1), (1,1,1), (1,0,1), (1,0,1)]
y = [1,1,0,0,1,1]
```

**Our X (predictors) data matrix.**

```
p = perceptron(X,y)
p.initialize()
p.train()
print p.predict((1,1,1))
print p.predict((1,0,1))
```

**We instantiate our perceptron class with the data from matrix X and vector y.**

**The weights for the predictors are initialized (as explained previously).**

**The perceptron model is trained. It will try to train until it either converges (no more errors) or it runs out of training runs (epochs).**

**We check what the perceptron would now predict given different values for the predictor variables. In the first case it will predict 0; in the second it predicts a 1.**

We'll zoom in on parts of the code that might not be so evident to grasp without further explanation. We'll start by explaining how the train\_observation() function works. This function has two large parts. The first is to calculate the prediction of an observation and compare it to the actual value. The second part is to change the weights if the prediction seems to be wrong.

```
def train_observation(self,X,y, error_count):
    result = np.dot(X, self.weights) > self.threshold
    error = y - result
    if error != 0:
        error_count += 1
        for index, value in enumerate(X):
            self.weights[index] += self.learning_rate * error * value
    return error_count
```

**The real value (y) is either 0 or 1; the prediction is also 0 or 1. If it's wrong we get an error of either 1 or -1.**

**The train observation function is run for every observation and will adjust the weights using the formula explained earlier.**

**A prediction is made for this observation. Because it's binary, this will be either 0 or 1.**

**In case we have a wrong prediction (an error), we need to adjust the model.**

**Adds 1 to error count.**

**Adjusts the weight for every predictor variable using the learning rate, the error, and the actual value of the predictor variable.**

**We return the error count because we need to evaluate it at the end of the epoch.**

**For every predictor variable in the input vector (X), we'll adjust its weight.**

The prediction (y) is calculated by multiplying the input vector of independent variables with their respective weights and summing up the terms (as in linear regression). Then this value is compared with the threshold. If it's larger than the threshold, the algorithm will give a 1 as output, and if it's less than the threshold, the algorithm gives 0 as output. Setting the threshold is a subjective thing and depends on your business case. Let's say you're predicting whether

someone has a certain lethal disease, with 1 being positive and 0 negative. In this case it's better to have a lower threshold: it's not as bad to be found positive and do a second investigation as it is to overlook the disease and let the patient die. The error is calculated, which will give the direction to the change of the weights.

---

```
result = np.dot(X, self.weights) > self.threshold
error = y - result
```

The weights are changed according to the sign of the error. The update is done with the learning rule for perceptrons. For every weight in the weight vector, you update its value with the following rule:

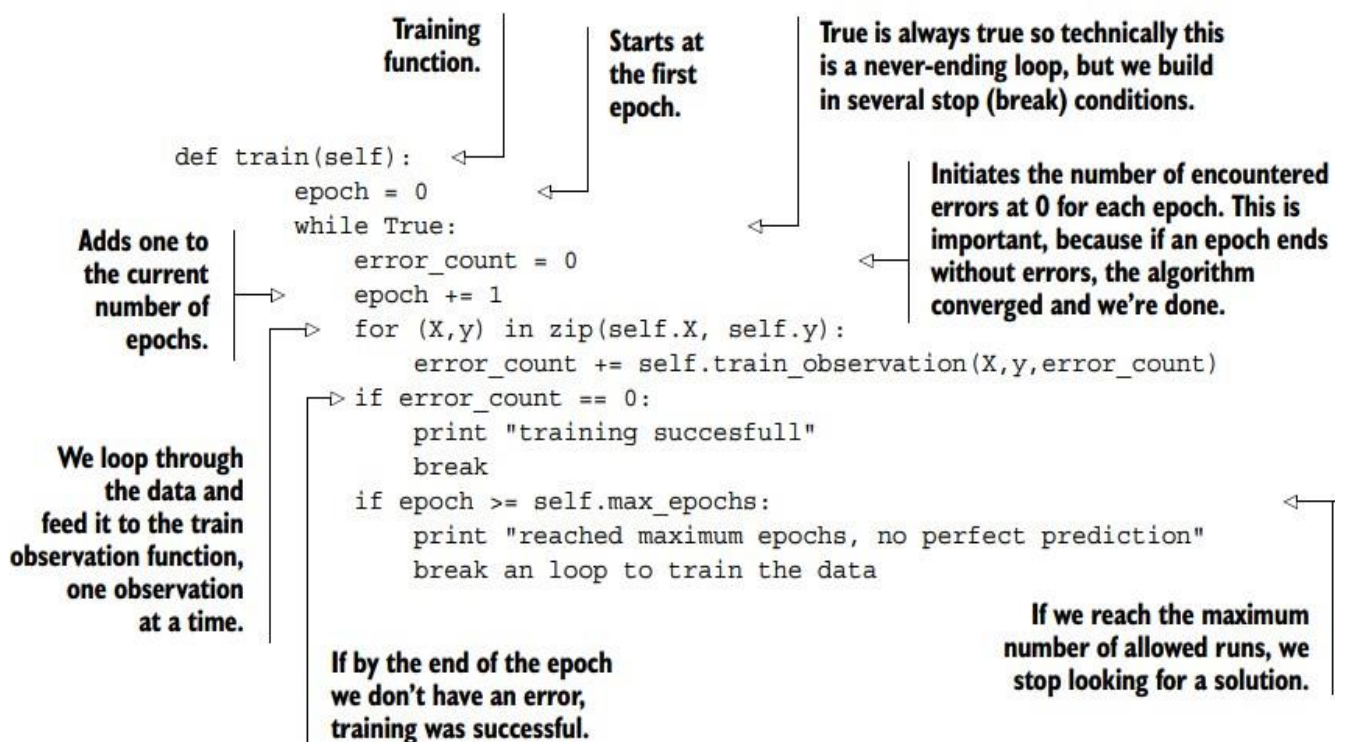
$$\Delta w_i = \alpha \varepsilon x_i$$

Where  $\Delta w_i$  is the amount that the weight needs to be changed,  $\alpha$  is the learning rate,  $\varepsilon$  is the error, and  $x_i$  is the  $i^{\text{th}}$  value in the input vector (the  $i^{\text{th}}$  predictor variable). The error count is a variable to keep track of how many observations are wrongly predicted in this epoch and is returned to the calling function. You add one observation to the error counter if the original prediction was wrong. An epoch is a single training run through all the observations.

```
if error != 0:
    error_count += 1
    for index, value in enumerate(X):
        self.weights[index] += self.learning_rate * error * value
```

The second function that we'll discuss in more detail is the `train()` function. This function has an internal loop that keeps on training the perceptron until it can either predict perfectly or until it has reached a certain number of training rounds (epochs), as shown in the following listing.

## Listing 4.2 Using train functions



Most online algorithms can also handle mini-batches; this way, you can feed them batches of 10 to 1,000 observations at once while using a sliding window to go over your data. You have three options:

- **Full batch learning (also called statistical learning)**—Feed the algorithm all the data at once.
- **Mini-batch learning**—Feed the algorithm a spoonful (100, 1000, ..., depending on what your hardware can handle) of observations at a time.
- **Online learning**—Feed the algorithm one observation at a time.

Online learning techniques are related to streaming algorithms, where you see every data point only once. Think about incoming Twitter data: it gets loaded into the algorithms, and then the observation (tweet) is discarded because the sheer number of incoming tweets of data might soon overwhelm the hardware. Online learning algorithms differ from streaming algorithms in that they can see the same observations multiple times. True, the online learning algorithms and streaming algorithms can both learn from observations one by one. Where they differ is that online algorithms are also used on a static data source as well as on a streaming data source by presenting the data in small batches (as small as a single observation), which enables you to go over the data multiple times. This isn't the case with a streaming algorithm, where data flows into the system and you need to do the calculations typically immediately. They're similar in that they handle only a few at a time.

## DIVIDING A LARGE MATRIX INTO MANY SMALL ONES:

By cutting a large data table into small matrices, for instance, we can still do a linear regression. The logic behind this matrix splitting and how a linear regression can be calculated with matrices can be found in the sidebar. It suffices to know for now that the Python libraries we're about to use will take care of the matrix splitting, and linear regression variable weights can be calculated using matrix calculus.

### Block matrices and matrix formula of linear regression coefficient estimation: -

Certain algorithms can be translated into algorithms that use blocks of matrices instead of full matrices. When you partition a matrix into a block matrix, you divide the full matrix into parts and work with the smaller parts instead of the full matrix. In this case you can load smaller matrices into memory and perform calculations, thereby avoiding an out-of-memory error. Figure 4.4 shows how you can rewrite matrix addition  $A + B$  into submatrices.

The formula in figure 4.4 shows that there's no difference between adding matrices  $A$  and  $B$  together in one step or first adding the upper half of the matrices and then adding the lower half.

All the common matrix and vector operations, such as multiplication, inversion, and singular value decomposition (a variable reduction technique like PCA), can be written in terms of block matrices.<sup>1</sup> Block matrix operations save memory by splitting the problem into smaller blocks and are easy to parallelize.

Although most numerical packages have highly optimized code, they work only with matrices that can fit into memory and will use block matrices in memory when advantageous. With out-of-memory matrices, they don't optimize this for you and it's up to you to partition the matrix into smaller matrices and to implement the block matrix version.

$$\begin{aligned}
 A + B &= \begin{bmatrix} a_{1,1} & \cdots & a_{1,m} \\ \vdots & \ddots & \vdots \\ a_{n,1} & \cdots & a_{n,m} \end{bmatrix} + \begin{bmatrix} b_{1,1} & \cdots & b_{1,m} \\ \vdots & \ddots & \vdots \\ b_{n,1} & \cdots & b_{n,m} \end{bmatrix} \\
 &= \begin{bmatrix} a_{1,1} & \cdots & a_{1,m} \\ \vdots & \ddots & \vdots \\ a_{j,1} & \cdots & a_{j,m} \\ \hline a_{j+1,1} & \cdots & a_{j+1,m} \\ \vdots & \ddots & \vdots \\ a_{n,1} & \cdots & a_{n,m} \end{bmatrix} + \begin{bmatrix} b_{1,1} & \cdots & b_{1,m} \\ \vdots & \ddots & \vdots \\ b_{j,1} & \cdots & b_{j,m} \\ \hline b_{j+1,1} & \cdots & b_{j+1,m} \\ \vdots & \ddots & \vdots \\ b_{n,1} & \cdots & b_{n,m} \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} + \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}
 \end{aligned}$$

**Figure 4.4** Block matrices can be used to calculate the sum of the matrices  $A$  and  $B$ . || Page

A linear regression is a way to predict continuous variables with a linear combination of its predictors; one of the most basic ways to perform the calculations is with a technique called ordinary least squares. The formula in matrix form is

$$\beta = (X^T X)^{-1} X^T y$$

where  $\beta$  is the coefficients you want to retrieve,  $X$  is the predictors, and  $y$  is the target variable.

The Python tools we have at our disposal to accomplish our task are the following:

- bcolz is a Python library that can store data arrays compactly and uses the hard drive when the array no longer fits into the main memory.
- Dask is a library that enables you to optimize the flow of calculations and makes performing calculations in parallel easier. It doesn't come packaged with the default Anaconda setup so make sure to use conda install dask on your virtual environment before running the code below. Note: some errors have been reported on importing Dask when using 64bit Python. Dask is dependent on a few other libraries (such as toolz), but the dependencies should be taken care of automatically by pip or conda.



The following listing demonstrates block matrix calculations with these libraries

**Listing 4.3 Block matrix calculations with bcolz and Dask libraries**

**Number of observations (scientific notation).  $1e4 = 10.000$ . Feel free to change this.**

```
import dask.array as da
import bcolz as bc
import numpy as np
import dask
```

```
n = 1e4
```

```
ar = bc.carray(np.arange(n).reshape(n/2,2), dtype='float64',
               rootdir = 'ar.bcolz', mode = 'w')
y = bc.carray(np.arange(n/2), dtype='float64', rootdir =
              'yy.bcolz', mode = 'w')
```

**Creates fake data: `np.arange(n).reshape(n/2,2)` creates a matrix of 5000 by 2 (because we set `n` to 10.000). `bc.carray = numpy` is an array extension that can swap to disc. This is also stored in a compressed way. `rootdir = 'ar.bcolz'` --> creates a file on disc in case out of RAM. You can check this on your file system next to this ipython file or whatever location you ran this code from. `mode = 'w'` --> is the write mode. `dtype = 'float64'` --> is the storage type of the data (which is float numbers).**

```
dax = da.from_array(ar, chunks=(5,5))
dy = da.from_array(y, chunks=(5,5))
```

**Block matrices are created for the predictor variables (`ar`) and target (`y`). A block matrix is a matrix cut in pieces (blocks). `da.from_array()` reads data from disc or RAM (wherever it resides currently). `chunks=(5,5)`: every block is a 5x5 matrix (unless < 5 observations or variables are left).**

**The `XTX` is defined (defining it as “lazy”) as the `X` matrix multiplied with its transposed version. This is a building block of the formula to do linear regression using matrix calculation.**

```
XTX = dax.T.dot(dax)
Xy = dax.T.dot(dy)
```

**`Xy` is the `y` vector multiplied with the transposed `X` matrix. Again the matrix is only defined, not calculated yet. This is also a building block of the formula to do linear regression using matrix calculation (see formula).**

```
coefficients = np.linalg.inv(XTX.compute()).dot(Xy.compute())
```

```
coef = da.from_array(coefficients, chunks=(5,5))
```

```
ar.flush()
y.flush()
```

**Flush memory data. It's no longer needed to have large matrices in memory.**

**The coefficients are also put into a block matrix. We got a numpy array back from the last step so we need to explicitly convert it back to a “da array.”**

```
predictions = dax.dot(coef).compute()
print predictions
```

**Score the model (make predictions).**

**The coefficients are calculated using the matrix linear regression function. `np.linalg.inv()` is the  $^{-1}$  in this function, or “inversion” of the matrix. `X.dot(y)` --> multiplies the matrix `X` with another matrix `y`.**

Note that you don't need to use a block matrix inversion because `XTX` is a square matrix with size `nr. of predictors * nr. of predictors`. This is fortunate because Dask doesn't yet support block matrix inversion.

## MAPREDUCE

MapReduce algorithms are easy to understand with an analogy: Imagine that you were asked to count all the votes for the national elections. Your country has 25 parties, 1,500 voting offices, and 2 million people. You could choose to gather all the voting tickets from every office individually and count them centrally, or you could ask the local offices to count the votes for the 25 parties and hand over the results to you, and you could then aggregate them by party.

Map reducers follow a similar process to the second way of working. They first map values to a key and then do an aggregation on that key during the reduce phase. Have a look at the following listing's pseudo code to get a better feeling for this.

#### Listing 4.4 MapReduce pseudo code example

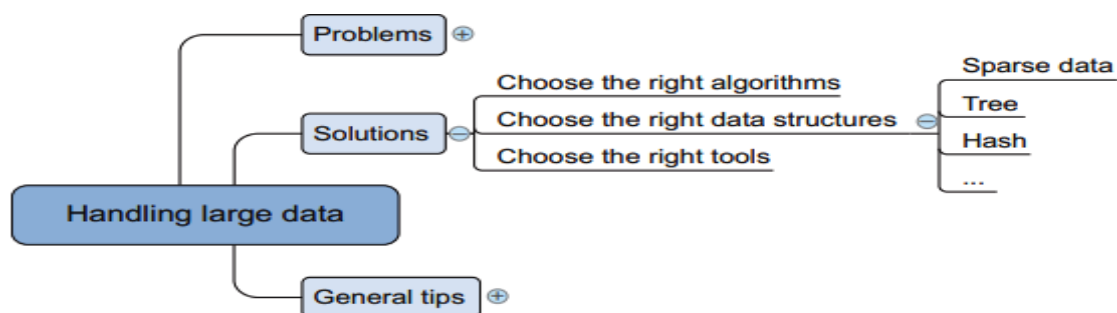
```
For each person in voting office:
    Yield (voted_party, 1)
For each vote in voting office:
    add_vote_to_party()
```

One of the advantages of MapReduce algorithms is that they're easy to parallelize and distribute. This explains their success in distributed environments such as Hadoop, but they can also be used on individual computers. When implementing MapReduce in Python, you don't need to start from scratch. A number of libraries have done most of the work for you, such as Hadoopy, Octopy, Disco, or Dumbo.

## 2. Choosing the right data structure: -

Algorithms can make or break your program, but the way you store your data is of equal importance. Data structures have different storage requirements, but also influence the performance of CRUD (create, read, update, and delete) and other operations on the data set.

Figure 4.5 shows you have many different data structures to choose from, three of which we'll discuss here: **sparse data**, **tree data**, and **hash data**. Let's first have a look at sparse data sets.



**Figure 4.5 Overview of data structures often applied in data science when working with large data**

## SPARSE DATA: -

A sparse data set contains relatively little information compared to its entries (observations). Look at figure 4.6: almost everything is "0" with just a single "1" present in the second observation on variable 9.

Data like this might look ridiculous, but this is often what you get when converting textual data to binary data. Imagine a set of 100,000 completely unrelated Twitter tweets. Most of them probably have fewer than 30 words, but together they might have hundreds or thousands of distinct words. In text mining we'll go through the process of cutting text documents into words

and storing them as vectors. But for now imagine what you'd get if every word was converted to a binary variable, with "1" representing "present in this tweet," and "0" meaning "not present in this tweet." This would result in sparse data indeed. The resulting large matrix can cause memory problems even though it contains little information.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

**Figure 4.6 Example of a sparse matrix: almost everything is 0; other values are the exception in a sparse matrix**

Luckily, data like this can be stored compacted. In the case of figure 4.6 it could look like this:

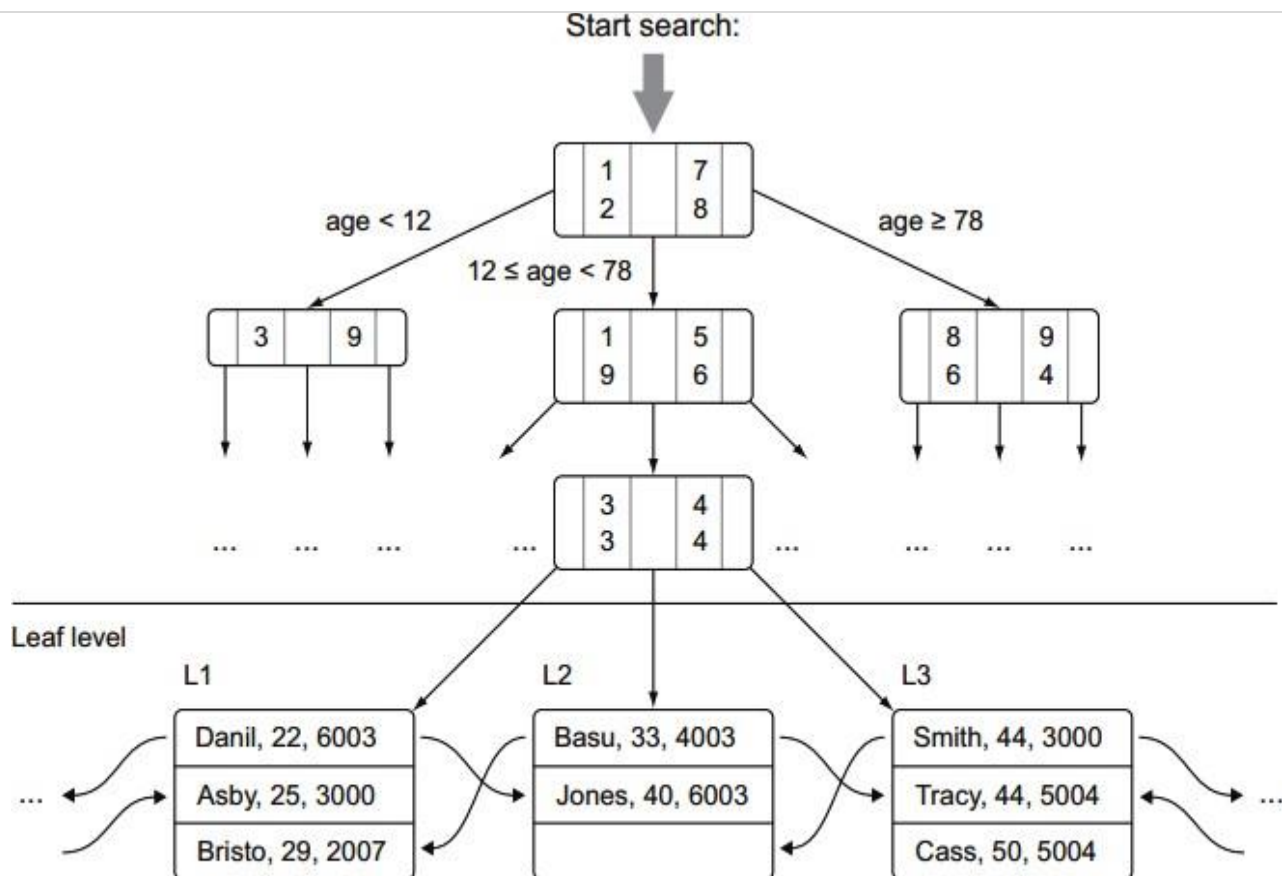
```
data = [(2,9,1)]
```

Row 2, column 9 holds the value 1.

Support for working with sparse matrices is growing in Python. Many algorithms now support or return sparse matrices.

## TREE STRUCTURES: -

Trees are a class of data structure that allows you to retrieve information much faster than scanning through a table. A tree always has a root value and sub trees of children, each with its children, and so on. Simple examples would be your own family tree or a biological tree and the way it splits into branches, twigs, and leaves. Simple decision rules make it easy to find the child tree in which your data resides. Look at figure 4.7 to see how a tree structure enables you to get to the relevant information quickly.



**Figure 4.7 Example of a tree data structure: decision rules such as age categories can be used to quickly locate a person in a family tree**

In figure 4.7 you start your search at the top and first choose an age category, because apparently that's the factor that cuts away the most alternatives. This goes on and on until you get what you're looking for.

Trees are also popular in databases. Databases prefer not to scan the table from the first line until the last, but to use a device called an index to avoid this. Indices are often based on data structures such as trees and hash tables to find observations faster. The use of an index speeds up the process of finding data enormously.

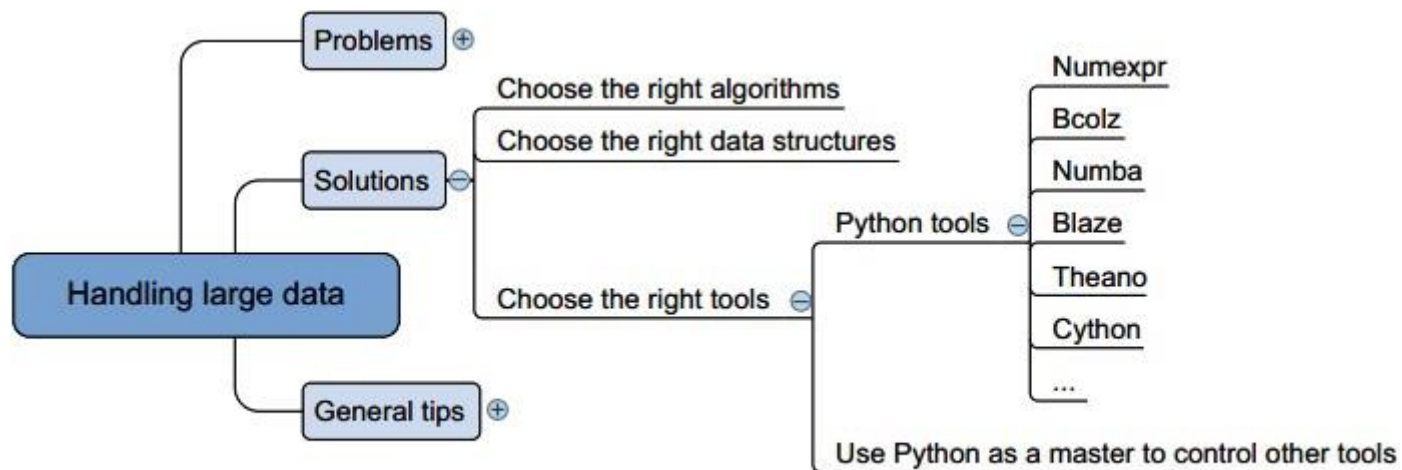
## HASH TABLES: -

Hash tables are data structures that calculate a key for every value in your data and put the keys in a bucket. This way you can quickly retrieve the information by looking in the right bucket when you encounter the data. Dictionaries in Python are a hash table implementation, and they're a close relative of key-value stores. Hash tables are used extensively in databases as indices for fast information retrieval.

### 3. Selecting the right tools: -

With the right class of algorithms and data structures in place, it's time to choose the right tool for the job. The right tool can be a Python library or at least a tool that's

controlled from Python, as shown figure 4.8. The number of helpful tools available is enormous, so we'll look at only a handful of them.



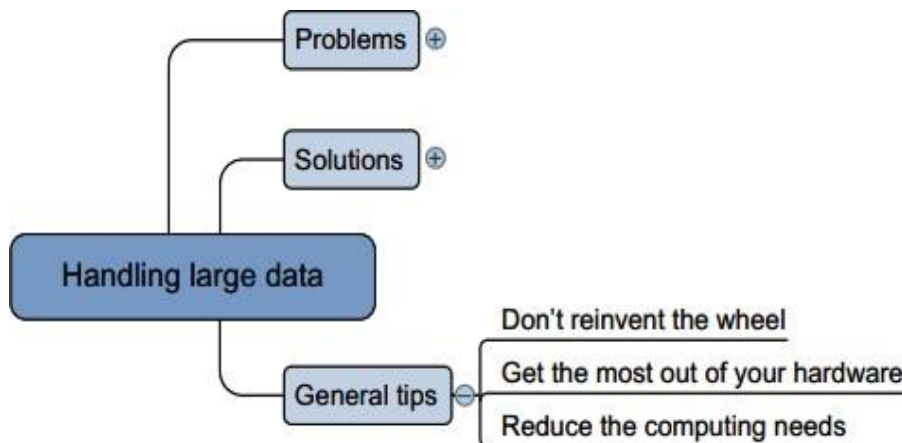
**Figure 4.8 Overview of tools that can be used when working with large data**

## General programming tips for dealing with large data sets: -

The tricks that work in a general programming context still apply for data science. Several might be worded slightly differently, but the principles are essentially the same for all programmers.

You can divide the general tricks into three parts, as shown in the figure 4.9 mind map:

- **Don't reinvent the wheel.** Use tools and libraries developed by others.
- **Get the most out of your hardware.** Your machine is never used to its full potential; with simple adaptations you can make it work harder.
- **Reduce the computing need.** Slim down your memory and processing needs as much as possible.



**Figure 4.9 Overview of general programming best practices when working with large data**

## 1. Don't reinvent the wheel: -

“Don't repeat anyone” is probably even better than “don't repeat yourself.” Add value with your actions: make sure that they matter. Solving a problem that has already been solved is a waste of time. As a data scientist, you have two large rules that can help you deal with large data and make you much more productive, to boot:

- **Exploit the power of databases.** The first reaction most data scientists have when working with large data sets is to prepare their analytical base tables inside a database. This method works well when the features you want to prepare are fairly simple. When this preparation involves advanced modeling, find out if it's possible to employ user- defined functions and procedures. The last example of this chapter is on integrating a database into your workflow.
- **Use optimized libraries.** Creating libraries like Mahout, Weka, and other machinelearning algorithms requires time and knowledge. They are highly optimized and incorporate best practices and state-of-the art technologies. Spend your time on getting things done, not on reinventing and repeating others people's efforts, unless it's for the sake of understanding how things work.

Then you must consider your hardware limitation.

## 2. Get the most out of your hardware: -

Resources on a computer can be idle, whereas other resources are over-utilized. This slows down programs and can even make them fail. Sometimes it's possible (and necessary) to shift the workload from an overtaxed resource to an underutilized resource using the following techniques:

- **Feed the CPU compressed data.** A simple trick to avoid CPU starvation is to feed the

CPU compressed data instead of the inflated (raw) data. This will shift more work from the hard disk to the CPU, which is exactly what you want to do, because a hard disk can't follow the CPU in most modern computer architectures.

- **Make use of the GPU.** Sometimes your CPU and not your memory is the bottleneck. If your computations are parallelizable, you can benefit from switching to the GPU. This has a much higher throughput for computations than a CPU. The GPU is enormously efficient in parallelizable jobs but has less cache than the CPU. But it's pointless to switch to the GPU when your hard disk is the problem. Several Python packages, such as Theano and NumbaPro, will use the GPU without much programming effort. If this doesn't suffice, you can use a CUDA (Compute Unified Device Architecture) package such as PyCUDA. It's also a well-known trick in bitcoin mining, if you're interested in creating your own money.
- **Use multiple threads.** It's still possible to parallelize computations on your CPU. You can achieve this with normal Python threads.

### 3. Reduce your computing needs: -

"Working smart + hard = achievement." This also applies to the programs you write. The best way to avoid having large data problems is by removing as much of the work as possible up front and letting the computer work only on the part that can't be skipped. The following list contains methods to help you achieve this:

- 
- **Profile your code and remediate slow pieces of code.** Not every piece of your code needs to be optimized; use a profiler to detect slow parts inside your program and remediate these parts.
  - **Use compiled code whenever possible, certainly when loops are involved.** Whenever possible use functions from packages that are optimized for numerical computations instead of implementing everything yourself. The code in these packages is often highly optimized and compiled.
  - **Otherwise, compile the code yourself.** If you can't use an existing package, use either a just-in-time compiler or implement the slowest parts of your code in a lower-level language such as C or Fortran and integrate this with your codebase. If you make the step to lower-level languages (languages that are closer to the universal computer bytecode), learn to work with computational libraries such as LAPACK, BLAST, Intel MKL, and ATLAS. These are highly optimized, and it's difficult to achieve similar performance to them.
  - **Avoid pulling data into memory.** When you work with data that doesn't fit in your memory, avoid pulling everything into memory. A simple way of doing this is by reading data in chunks and parsing the data on the fly. This won't work on every algorithm but enables calculations on extremely large data sets.
  - **Use generators to avoid intermediate data storage.** Generators help you return data



per observation instead of in batches. This way you avoid storing intermediate results.

- **Use as little data as possible.** If no large-scale algorithm is available and you aren't willing to implement such a technique yourself, then you can still train your data on only a sample of the original data.
- **Use your math skills to simplify calculations as much as possible.** Take the following equation, for example:  $(a + b)^2 = a^2 + 2ab + b^2$ . The left side will be computed much faster than the right side of the equation; even for this trivial example, it could make a difference when talking about big chunks of data.

#### UNIT WISE IMPORTANT QUESTIONS: -

1. Explain the problems you face when handling large data.
  2. Explain the solutions for handling large data sets.
  3. Discuss in detail about block matrices
  4. What is Map Reduce? Explain in detail.
  5. Explain in details about online learning algorithms
  6. What are problems encountered when working with more data than can fit in memory? Explain
  7. How to choose the right data structure? Explain
  8. List and explain different libraries that can help you deal with large data.
  9. "Don't reinvent the wheel". Justify the statement
  10. Explain about general programming best practices when working with large data in detail.
  11. What are the different techniques to adapt algorithms to large data sets? Explain.
- 

### CASE STUDIES ON DS PROJECTS FOR PREDICTING MALICIOUS URLS FOR BUILDING RECOMMENDER SYSTEMS

Building recommender systems for predicting malicious URLs involves leveraging data science techniques to provide effective recommendations or alerts based on the characteristics of URLs. Here are a few hypothetical case studies to illustrate different approaches:

#### Case Study 1: Content-Based Recommender System

**Problem Statement:** Develop a content-based recommender system to suggest whether a URL is likely to be malicious based on its content and structural features.

**Dataset:** Use a dataset of URLs labeled as malicious or benign, along with features extracted from the URLs such as domain age, URL length, presence of special characters, etc.

**Approach:**

1. **Feature Extraction:** Extract relevant features from the URLs, focusing on both structural (e.g., domain age) and content-based features (e.g., lexical analysis of URL path).
2. **Content Representation:** Represent URLs using embeddings or vector representations (e.g., TF-IDF, Word2Vec) to capture semantic meaning and structural characteristics.
3. **Similarity Calculation:** Calculate similarity between the new URL and known malicious URLs using cosine similarity or other distance metrics.



4. **Thresholding:** Define a threshold based on similarity scores to classify URLs as potentially malicious or benign.
5. **Evaluation:** Evaluate the recommender system's performance using metrics such as precision, recall, and F1-score. Consider using techniques like cross-validation to validate the model.
6. **Deployment:** Deploy the content-based recommender system in a real-time environment where it can provide recommendations or alerts based on new URLs encountered.

## Case Study 2: Collaborative Filtering Approach

**Problem Statement:** Develop a collaborative filtering recommender system to predict whether a URL is malicious based on patterns learned from user behaviors or expert feedback.

**Dataset:** Utilize a dataset containing user feedback or expert ratings on URLs, indicating whether they are malicious or benign.

### Approach:

1. **User-Item Matrix:** Create a user-item matrix where rows represent users (or experts) and columns represent URLs, with values indicating feedback (e.g., ratings, binary labels).
2. **Model Selection:** Choose a collaborative filtering technique such as matrix factorization (e.g., Singular Value Decomposition, Alternating Least Squares) or neighborhood-based methods (e.g., k-Nearest Neighbors).
3. **Training:** Train the collaborative filtering model on the user-item matrix to learn latent factors representing user preferences or expert judgments regarding malicious URLs.
4. **Prediction:** Predict the likelihood of a new URL being malicious based on the learned patterns from the user-item matrix.
5. **Evaluation:** Evaluate the recommender system's performance using standard evaluation metrics for recommendation systems, such as Mean Absolute Error (MAE) or Root Mean Squared Error (RMSE).
6. **Deployment:** Deploy the collaborative filtering recommender system in a production environment, ensuring it can handle real-time recommendation tasks and provide reliable predictions.

## Case Study 3: Hybrid Approach

**Problem Statement:** Develop a hybrid recommender system combining content-based and collaborative filtering approaches for predicting malicious URLs.

**Dataset:** Use a comprehensive dataset combining structural features of URLs and user feedback or expert ratings.

### Approach:

1. **Feature Engineering:** Extract a combination of features including structural (e.g., domain age, URL length) and content-based features (e.g., semantic embeddings of URL content).
2. **Content-Based Filtering:** Develop a content-based filtering component to recommend whether a URL is malicious based on its features and similarity to known malicious URLs.
3. **Collaborative Filtering:** Develop a collaborative filtering component to leverage user feedback or expert ratings to predict the maliciousness of URLs.
4. **Hybrid Model Fusion:** Combine predictions from both content-based and collaborative filtering components using ensemble techniques (e.g., weighted average, stacking) to produce a final recommendation.
5. **Evaluation:** Evaluate the hybrid recommender system's performance using appropriate metrics, comparing it with individual components for validation and improvement.
6. **Deployment:** Deploy the hybrid recommender system in a production environment, ensuring scalability and real-time performance for providing recommendations or alerts on malicious URLs.

