

Module 5

Chapter 5 (T2) – Spark and Big Data Analytics

Prepared By:

Mr. Rajesh Nayak

Department of Artificial Intelligence and Data Science

• **5.1 Introduction**

- Pig or Hive are high-level scripting languages used with Apache Hadoop. They have SQL like commands for queries.
- These commands are translated to Map and Reduce tasks before execution.
- They support queries, built-in functions, UDFs, and run ETL processes.
- Berkeley's Algorithms, Machines and Peoples Laboratory (AMP) developed Berkeley Data Analytics Stack (BDAS) which supports efficient, large-scale in-memory data processing, and includes applications fulfilling three fundamental processing requirements: accuracy, time and cost.
- AMP first developed Spark in 2009 and later passed on the project to Apache. A new version is Spark 2.3.1. (as of present market Spark 3.5.5).

- **5.2 Spark**

- Apache Spark is a fast and general compute engine. It powers the analytics application up to 100 times faster.
- It supports HDFS compatible data and has simple and expressive programming model.
- Expressive program model implements a number of mathematical and logic operations with smaller and easier written codes which a compiler as well as programmer can easily understand.
- Spark runs on both Windows and Unix-like systems, such as Linux and Mac OS. Java is essential for running Spark applications.

- **5.2.1 Introduction to Big Data Tool - Spark**
- Following figure 5.1 shows main components in Apache Software Foundation's Spark framework.

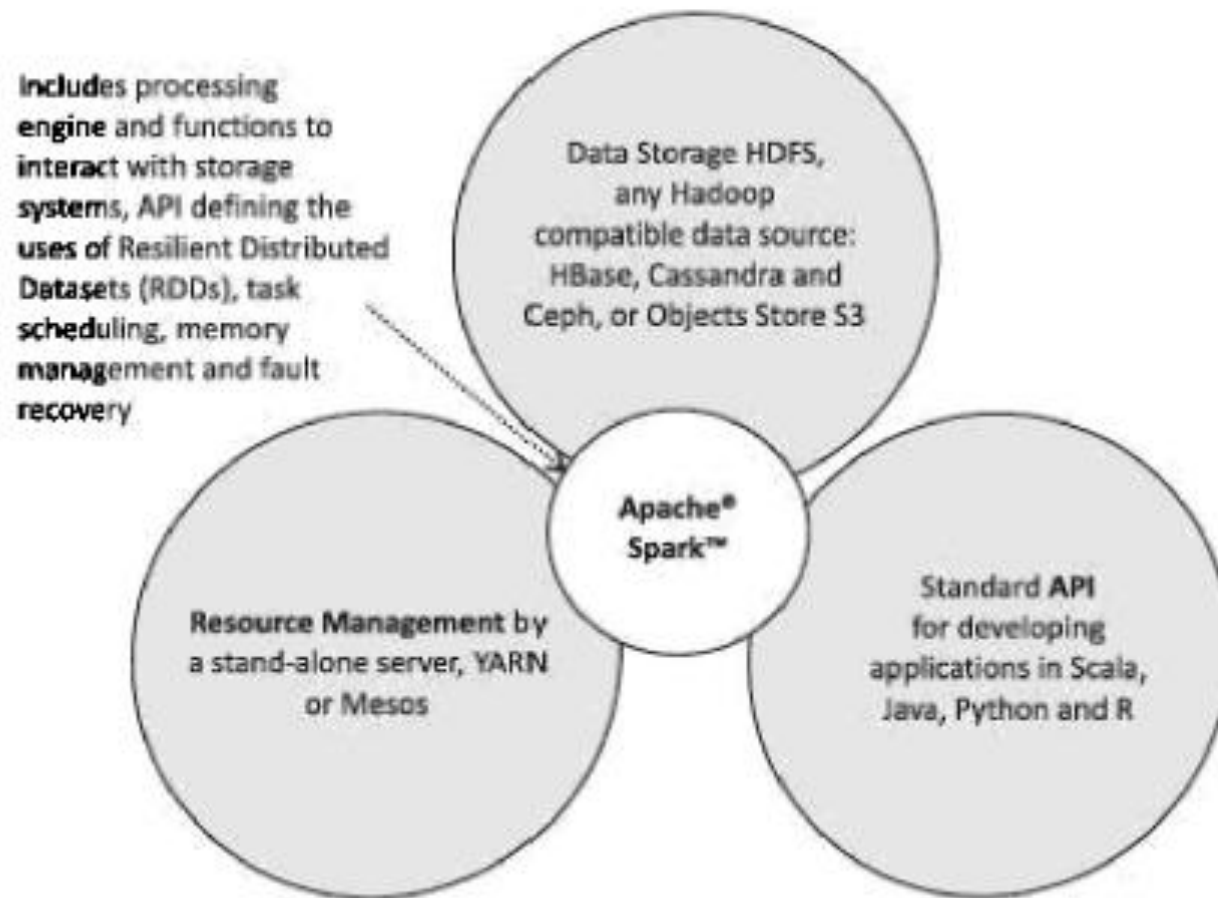


Figure 5.1 Main components of the Spark architecture

- **5.2.1 Introduction to Big Data Tool - Spark**
- Main components of Spark architecture are:
 1. Spark HDFS file system for data storage: Storage is at an HDFS or Hadoop compatible data source (such as HDFS, Hbase, Cassandra, Ceph), or at the Object Store S3.
 2. Spark standard API enables the creation of applications using Scala, Java, Python and R.
 3. Spark resource management can be at a stand-alone server or it can be on a distributed computing framework, such as YARN or Mesos.
- Ceph refers to an open source, scalable object, block, file storing unified system. It provides for dynamic replication and redistribution.
- Apache Mesos is an open source project developed at University of Berkeley, and now passed to Apache. It manages computing clusters. Its implementation is in C++. Mesos provides resources. Each resource contains list of agents. Each agent has the Hadoop and MapReduce executors.
- S3 (Simple Storage Service) refers to an Amazon Web Service on the cloud named S3 which provides Object Stores for data.

- **5.2.1.1 Features of Spark**
- Main features of Spark are shown in figure.

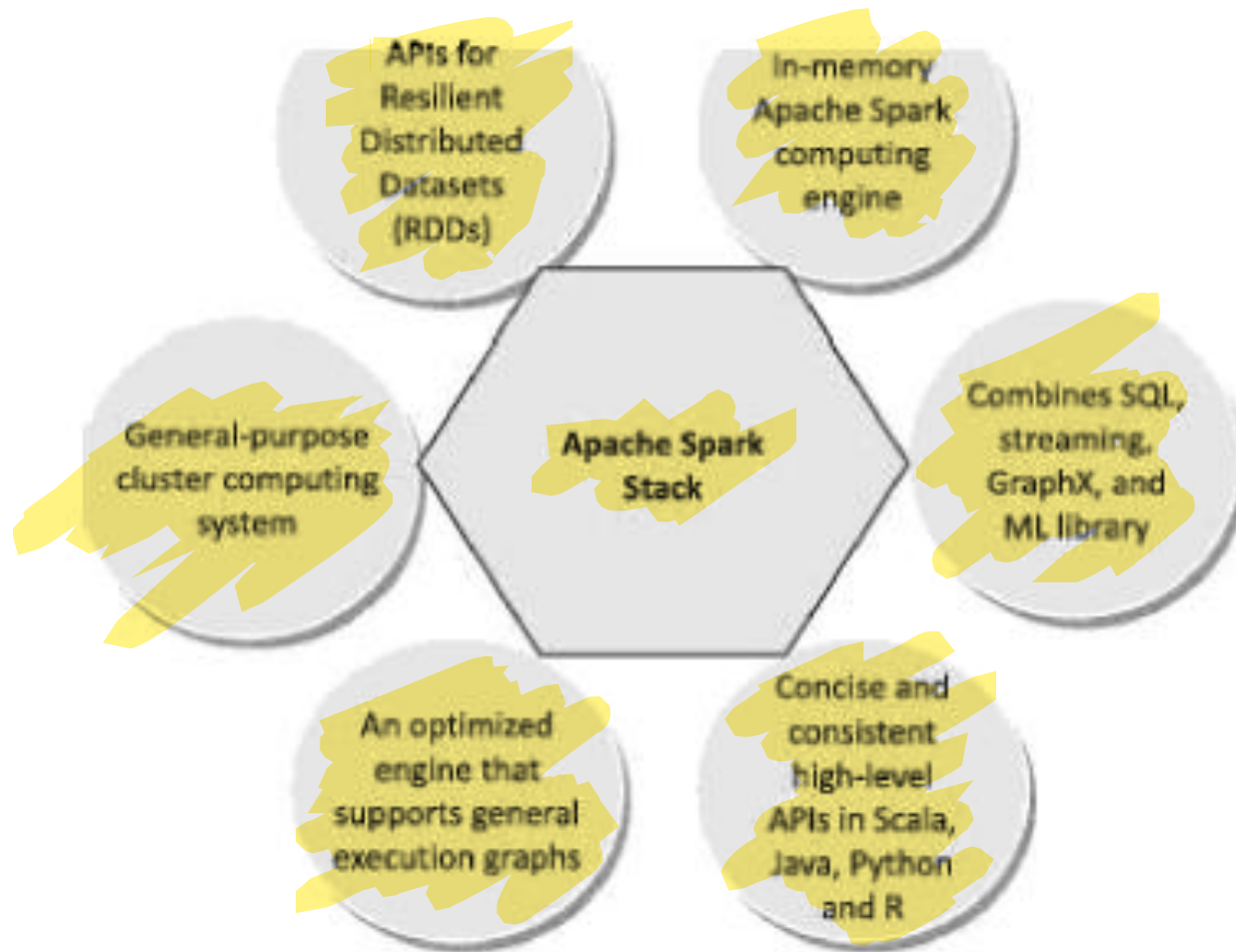


Figure 5.2 Main features of Spark

- **5.2.1.1 Features of Spark**

- **Main features of Spark are:**

1. Spark provisions for creating applications that use the complex data. In-memory Apache Spark computing engine enables up to 100 times performance with respect to Hadoop.
2. Execution engine uses both in-memory and on-disk computing. Intermediate results save in-memory and spill over to disk.
3. Data uploading from an Object Store for immediate use as a Spark object instance. Spark service interface sets up the Object Store.
4. Provides high performance when an application accesses memory cache from the disk.
5. Contain API to define Resilient Distributed Datasets (RDDs). RDD is a programming abstraction. It represents collection of Object Stores distributed across many compute nodes for parallel processing. Data is stored in different partitions and RDD is fault tolerant.
6. Processes any kind of data arriving from various sources.
7. Supports many new functions in addition to Map and Reduce functions.
8. Optimizes data processing performance by slowing the evaluation of Big Data queries.
9. Provides concise and consistent APIs in Scala, Java and Python. Spark codes are in Scala and run on JVM environment.
10. Supports Scala, Java, Python, Clojure and R languages.
11. Provides powerful tool to analyze data interactively using shell which is available in either Scala or Python.

- **5.2.1.2 Spark Software Stack**
- Figure 5.3 shows 5-layer architecture of running applications when using Spark stack.

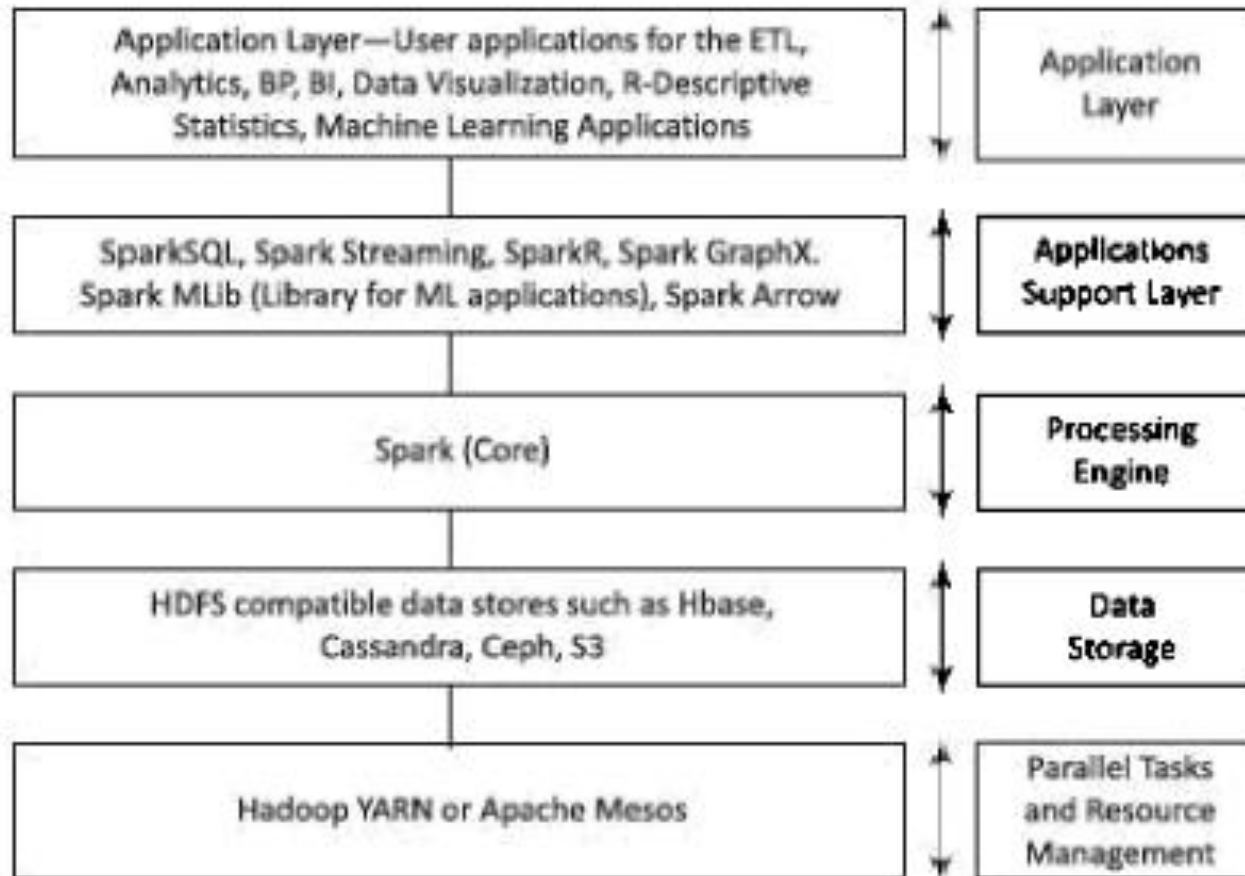


Figure 5.3 Five-layer architecture for running applications using Spark stack

- **5.2.1.2 Spark Software Stack**

The main components of Spark stack are SQL, Streaming, R, GraphX, MLib and Arrow at the applications support layer. Spark core is the processing engine. Data Store provides the data to the processing engine. Hadoop, YARN or Mesos facilitates the parallel running of the tasks and the management and scheduling of the resources.

Spark Stack

Spark stack imbibes generality to Spark. Grouping of the following forms Spark stack:

Spark SQL for the structured data. The SQL runs the queries on Spark data in the traditional business analytics and visualization applications. Spark SQL enables Spark datasets to use JDBC or ODBC API. HQL queries also run in Spark SQL. Runs UDFs for inline SQL, distributed DataFrames, Parquet, Hive and Cassandra Data Stores.

Spark Streaming is for processing real-time streaming data. Processing is based on micro-batches style of computing and processing. Streaming uses the DStream which is basically a series of RDDs, to process the real-time data.

- **5.2.1.2 Spark Software Stack**

SparkR is an R package used as light-weight front end for Apache Spark from R. Spark API uses SparkR through the RDD class. A user can interactively run the jobs from the R shell on a cluster. An RDD API is in the distributed lists in R.

Spark MLib is Spark's scalable machine learning library. It consists of common learning algorithms and utilities. MLib includes classification, regression, clustering, collaborative filtering, dimensionality reduction and optimization primitives. MLib applies in recommendation systems, clustering and classification using Spark.

Spark GraphX is an API for graphs. GraphX extends the Spark RDD by introducing the Resilient Distributed Property. GraphX computations use fundamental operators (e.g., subgraph, joinVertices and aggregateMessages). GraphX uses a collection of graph algorithms for programming. Graph analytics tasks are created with ease using GraphX.

Spark Arrow for columnar in-memory analytics and enabling usages of vectorised UDFs (VUDFs). The Arrow enables high performance Python UDFs for SerDe and data pipelines.

- **5.3 Introduction to Data Analysis with Spark**
- “Analysis of data is a process of inspecting, cleaning, transforming and modelling data with the goal of discovering useful information, suggesting conclusions and supporting decision-making.”

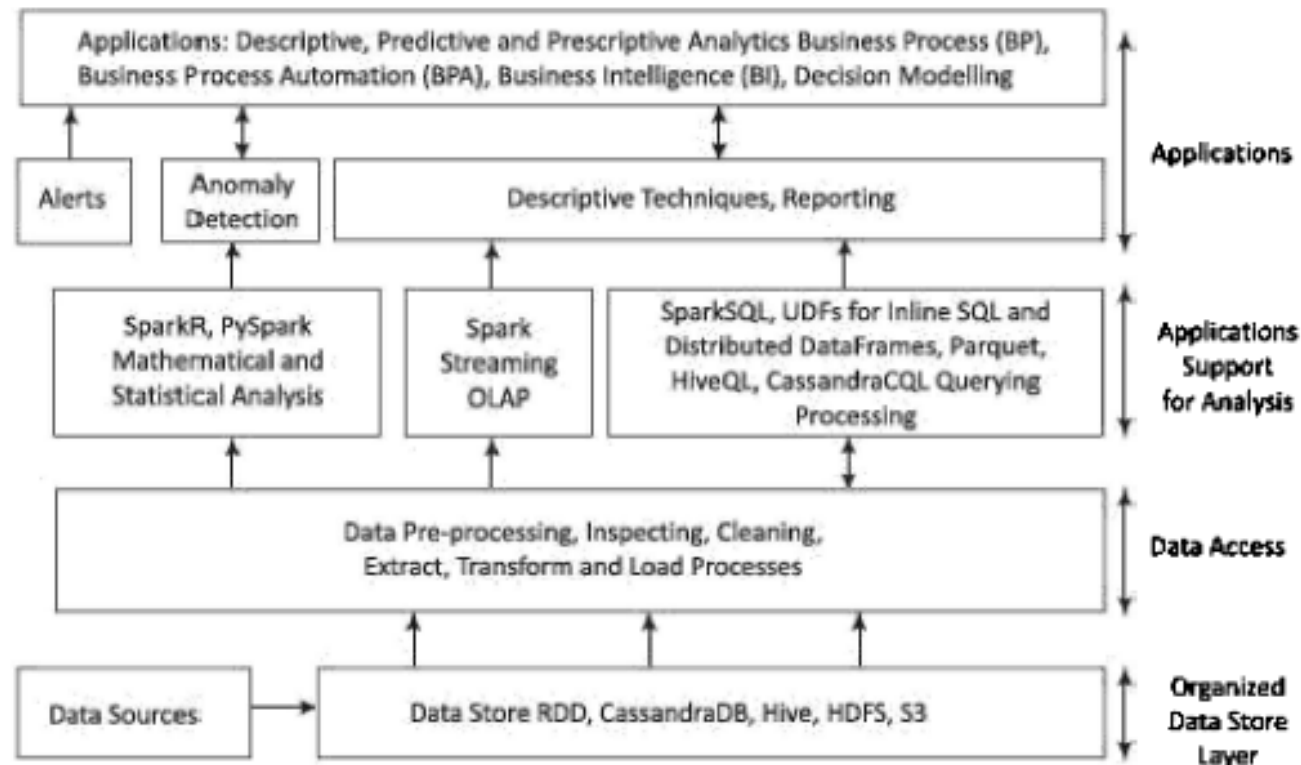


Figure 5.4 Steps between acquisition of data from different sources and its applications

• 5.3 Introduction to Data Analysis with Spark

Following are the steps for analyzing the data:

1. **Data Storage:** Store of data from the multiple sources after acquisition. The Big Data storage may be in HDFS compatible files, Cassandra, Hive, HDFS or S3.
2. **Data pre-processing:** This step requires:
 - (a) dropping out of range, inconsistent and outlier values,
 - (b) filtering unreliable, irrelevant and redundant information,
 - (c) data cleaning, editing, reduction and/or wrangling,
 - (d) data-validation, transformation or transcoding.
3. **Extract, transform and Load (ETL))** for the analysis
4. **Mathematical and statistical analysis** of the data obtained after querying relevant data needing the analysis, or OLAP,
5. **Applications of analyzed data**, for example, descriptive, predictive and prescriptive analytics, business processes (BPs), business process automation (BPA), business intelligence (BI), decision modelling and knowledge discovery.

• 5.3.1 Spark SQL

Spark SQL is a component of Spark Big Data Stack. Spark SQL components are DataFrames (SchemaRDDs), SQLContext and JDBC server. Spark SQL at Spark does the following:

1. Runs SQL like scripts for query processing, using *catalyst optimizer* and *tungsten execution engine*
2. Processes structured data
3. Provides flexible APIs for support for many types of data sources
4. ETL operations by creating ETL pipeline on the data from different file-formats, such as JSON, Parquet, Hive, Cassandra and then run ad-hoc querying.

• 5.3.1 Spark SQL

Spark SQL has the following features for analysis:

1. SparkR, PySpark, Python, Java and other language support for coding for data analysis.
2. Provisioning of JDBC and ODBC APIs: Applications in Java and Microsoft programs (such as Excel) need to connect to databases using JDBC (Object Database Connectivity) and ODBC (Object Database Connectivity). Spark APIs enable that connectivity,
3. Spark SQL enables users to extract their data from different formats, such as Hive, JSON and Parquet, and then transform that into required formats for ad hoc querying. [Ad hoc query is a query 'just for this purpose' or query 'on the fly.' For example, `var newToyQuery = "SELECT * FROM table WHERE id = " + toy_puzzleId`. The result will be different each time this code executes, depending on the value of `toy_puzzleId`.
4. Spark SQL processing support inclusion of Hive. Hive support enables the use of Hive tables, database and data warehouse, UDFs and SerDe (serialization and deserialization).
5. Spark SQL supports HiveQL and Cassandra CQL for query processing.
6. Spark Streaming for support to OLTP and structured streaming.

• 5.3.1 Spark SQL

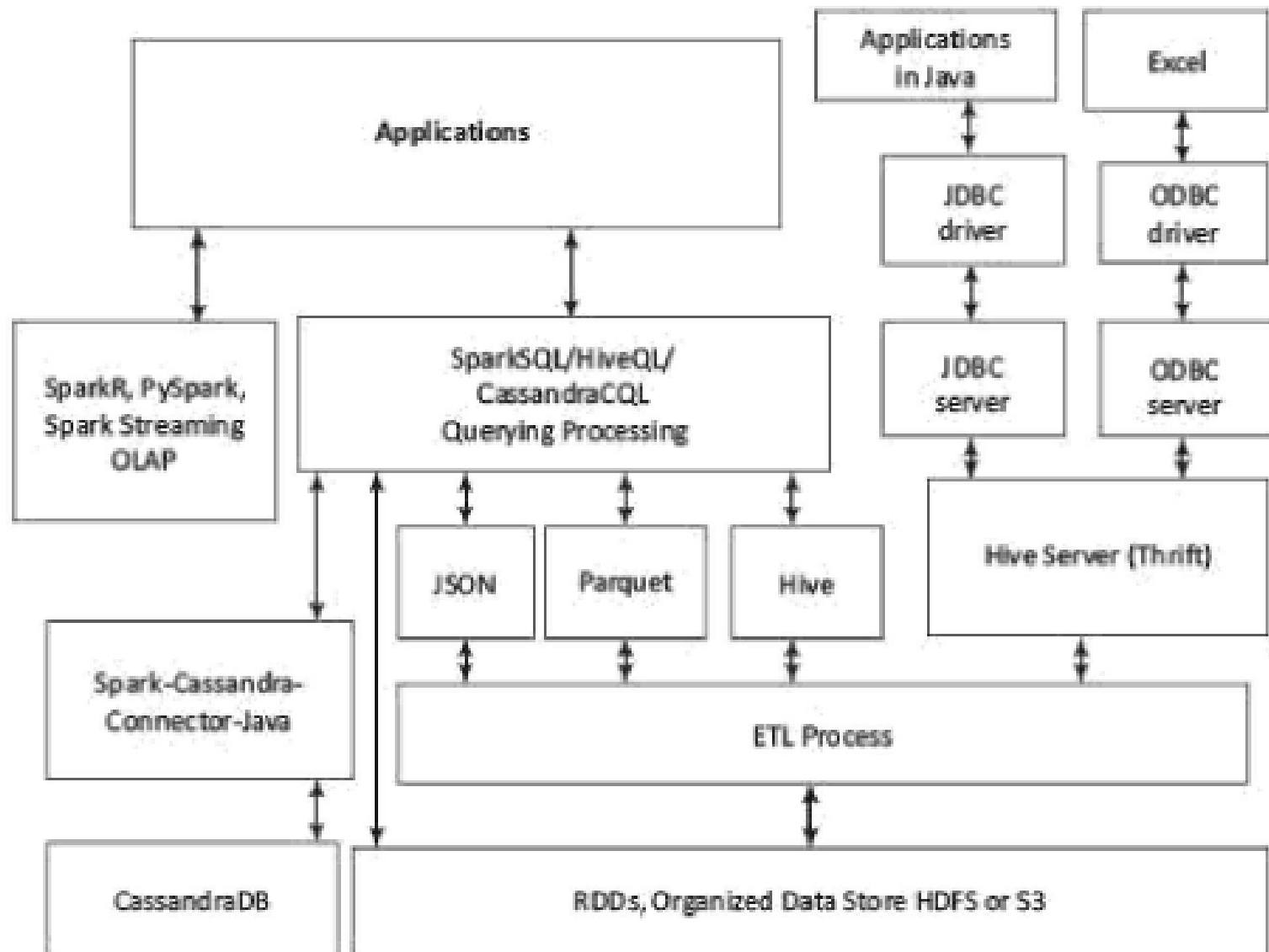


Figure 5.5 Connectivity between the applications and Spark SQL

- **5.3.2 Using Python Advanced Features with Spark SQL**

Python is a general purpose, interpreted, interactive, object oriented and high level programming language. Python defines the basic data types, containers, lists, dictionaries, sets, tuples, functions and classes. Python Standard Library is very extensive. The libraries for regular expressions, documentation generation, unit testing, web browsers, threading, databases, CGI, email, image manipulation and a lot of other functionalities are available in Python.

Python programming is a strong combination of performance and features in the same bundle of codes. Spark SQL binds with Python easily. Python has the expressive program statements. Spark SQL features together with Python help a programmer to build challenging applications for Big Data.

5.3.2.1 Python Libraries for Analysis

NumPy and SciPy are open source downloadable libraries for numerical (Num) analysis and scientific (Sci) computations in Python (Py). Python has open source library packages, NumPy, SciPy, Scikit-learn, Pandas and StatsModel, which are widely used for data analysis. Python library, matplotlib functions plot the mathematical functions.

Spark added a Python API support for UDFs. The functions take one row at a time. That requires overhead (additional codes) for SerDe. Earlier data pipelines first defined the UDFs in Java or Scala, and then invoked them from Python. Spark 2.3 provisions for vectorized UDFs (VUDFs) and Apache Arrow facilitates VUDFs, which enables high performance Python UDFs for SerDe and data pipelines.

NumPy NumPy includes (i) N-dimensional array object, array and vector mathematics; (ii) linear algebraic functions, Fourier transform and random number functions; (iii) sophisticated (broadcasting) functions; and (iv) tools for integrating with C/C++ and Fortran codes.

NumPy provides multi-dimensional efficient containers of generic data and definitions of arbitrary data types. NumPy integrates easily with a wide variety of databases. NumPy provides import, export (load/save) files, creation of arrays, inspection of properties, copying, sorting and reshaping, addition and removal of elements in the arrays, indexing, sub-setting and slicing of the arrays, scalar and vector mathematics (such as +, -, ×, ÷, power, sqr, sin, log, ceil – round up to nearest int, floor – round down up to the nearest int, round – round to nearest integer). NumPy also provides statistical functions.

Table 5.1 gives the examples of NumPy functions for data analysis problems.

Table 5.1 Examples of NumPy functions for data analysis problems

Function	Description	Function	Description
<code>np.loadtxt('file.txt')</code>	Loads a text file	<code>np.mean(arr, axis = 0)</code>	Returns mean along a specific axis
<code>np.genfromtxt('file.csv', delimiter= ',')</code>	Loads a csv file with comma as the delimiter between records	<code>np.sum(); np.min(),</code>	Returns the sum and minimum of the array
<code>np.savetxt('file.txt', arr, delimiter= ' ')</code>	Saves a text file which is an array of strings separated by a space each.	<code>np.max(arr, axis = 0)</code>	Returns the maximum along a specific axis
<code>np.genfromtxt('file.csv', arr, delimiter= ',')</code>	Saves a CSV file which is an array of strings separated by a comma each.	<code>np.var(arr)</code>	Returns the variance of array
<code>arr.sort()</code>	Sorts an array	<code>np.std(arr, axis = 0)</code>	Returns the standard deviation of a specific axis
<code>np.add (arr1, arr2)</code>	Performs a vector addition of array 1 and array 2	<code>np.corr()</code>	Returns the correlation coefficient

SciPy SciPy adds on top of NumPy. It includes MATLAB files and special functions, such as routines for numerical integration and optimization. SciPy defines some useful functions for computing distances between a set of points.

SciPy includes (i) interactions with NumPy, (ii) creation of dense and open mesh grids, (iii) shape manipulation functions, (iv) polynomial and vectoring functions, (v) real and imaginary functions, and casting an object to a data type, and (vi) matrix creation and matrices routines and usages of sparse matrices.

Table 5.2 gives few examples of SciPy functions for scientific computational problems.

Table 5.2 Examples of SciPy functions for data analysis problems

Function	Description	Function	Description
<code>np.c_[b, c]</code>	Create Stacked column-wise array	<code>np.cast ['f'] (np.pi)</code>	Casts an object into a data type
<code>b.flatten()</code>	Flattens the array	<code>from numpy import polyID p = polyID ([2, 3, 4])</code>	Creates a polynomial object p
<code>np.vsplit (c, 2)</code> and <code>np.hsplit (d, 2)</code>	Functions for vertically splitting and horizontally splitting the array at the end of the second index	<code>A.I, A.T, A.H</code>	Inverses, transposes, and conjugate transposes the matrix, A
<code>linalg.det(A)</code>	Returns the determinate of A	<code>np.select ([c<4], [c*2])</code>	Returns values from a list of arrays depending on the conditions
<code>np.img(c)</code>	Returns the imaginary part of the array elements	<code>A = np.matrix ([3, 4], [5, 6])</code>	Creates a matrix

Panda Panda derives its name from usages of a data structure called Panel. The first three characters *Pan* in Panda stand for the term 'panel'. The next two characters *da* in Panda stand for data. The Panda package considers three data structures: Series, DataFrame and Panel.

DataFrame is a container for Series. Panel is a container for DataFrame objects. The Panel objects can be inserted or removed similar to as in a dictionary. DataFrame may be a DataFrame of statistical or observed datasets. The data need not be labeled. This means datasets, objects and DataFrames can be placed into a Panda data structure without labels.

Panel is a container for three-dimensional data. Panel is a widely used term in econometrics. Three axes describe operations involving panel data. For example, panel data in econometric analysis. Items can be considered as along axis 0 of an inside DataFrame (set of columns). Index (rows) of each of the DataFrames correspond to axis 1 (major axis). Columns of each of the DataFrames correspond to axis 2 (minor axis). Panel 4D consists of labels as 0 - axis, items - axis 1, major axis - axis 2 and minor axis - axis 3.

Figure 5.7 shows the main features of Pandas package.



Figure 5.7 Main features of Panda for data analysis

Pandas package includes the following provisions:

1. Database style DataFrames merge, join, and concatenation of objects
2. RPy interface for R functions plus additional functions
3. Panda ecosystem has statistics, machine learning, integrated development environment (IDE), API and several out of core features
4. SQL like features: SELECT, WHERE, GROUPBY, JOIN, UNION, UPDATE and DELETE
5. GroupBy feature of split-apply-combine with the steps as: (i) an object such as table, file or document splits into groups, (ii) iterate through the groups and select a group for aggregation, transformation and/or filtration. The instance method can be dispatched and applied in a manner similar to the aggregation/transformation function
6. Size mutability, which means that columns can be inserted and deleted from DataFrame and higher dimensional objects
7. Slicing and dicing a collection of DataFrame objects. The names of axes can be somewhat arbitrary in a Panel. (Arbitrary means the axes need not be named a1, a2, ..., an, and can be year, car model, sales, ...). If slicing function slices the first dimension, the lower dimension objects are obtained.

5.3.2.2 User-Defined Functions (UDFs)

The functions take one row at a time. This requires overhead for SerDe. Data exchanges take place between Python and JVM. Earlier the data pipeline (between data and application) defined the UDFs in Java or Scala, and then invoked them from Python while using Python libraries for analysis or other application. SparkSQL UDFs enable registering of themselves in Python, Java and Scala.

The SQL calls the UDFs. This is a very popular way to expose advanced functionality to SQL users. User codes call the registered UDFs into the SQL statements without writing the detailed codes.

5.3.2.3 Vectorized User Defined Functions (VUDFs)

Python UDFs express data in detail. Therefore, Python UDFs, block-level UDFs with block-level arguments and return types, conversions or transformations are widely used in ETL or ML applications. Spark Arrow facilitates columnar in-memory analytics, which results in high performance of Python UDFs, SerDe and data pipelines.

VUDFs use series data structure (meaning one-dimensional array or tuples). Spark 2.3 (2018) provisions for using vectorized UDFs (VUDFs). Apache Arrow 0.8.0 (release date December 18, 2017) facilitates usages of VUDFs. Pandas UDF, *pandas_UDF* uses the function to create a VUDF with (i) pandas.Series as input to the UDF, (ii) pandas.Series as output from the UDF, (iii) no grouping using GroupBy, (iv) output size same as input, and (v) returns the same data types as specified type in return pandas.Series.

5.3.2.4 Grouped Vectorized UDFs (GVUDFs)

Grouped Vectorized UDFs (GVUDFs) use Panda library *split-apply-combine* pattern in data analysis. The GVUDF group function operates on all the data for a group, such as operate on all the data, “for each car showroom, compute yearly sales”.

GVUDF steps are:

1. Splits a Spark DataFrame into groups based on the conditions specified in the groupby operator
2. Applies a vectorized user-defined function (`pandas.DataFrame -> pandas.DataFrame`) to each group
3. Combines into new group
4. Returns the results as a new Spark DataFrame

Pandas GVUDF, `pandas_GVUDF`, (i) uses the function similar to `pandas_VUDF`, (ii) `pandas.DataFrame` as input to the GVUDF, (ii) `pandas.DataFrame` as output from the GVUDF, (iii) grouping semantics defined using clause `GroupBy`, (iv) output size can be any and can be grouped and can be distinct from input, and (v) returns data type is a `StructType`. The type defines a column name and the type of the returned `pandas.DataFrame`.

5.3.3 Data Analysis Operations

Examples of operations required in the above analysis are given below:

1. Filtering single and multiple columns
2. Creating a top-ten list with values or percentages
3. Setting up sub-totals
4. Creating multiple-field criteria filters
5. Creating unique lists from repeating field data
6. Finding duplicate data with specialized arrays, and using the remove duplicates command and removing outliers
7. Multiple key sorting
8. Counting the number of unique items in a list
9. Using SUMIF and COUNTIF functions
10. Working with database functions, such as DSUM and DMAX
11. Converting lists to tables.

5.3.3.1 Removing Outliers for Data Quality Improvement for Analysis

Outliers are data which appear as they do not belong to the data set. The Outliers are generally results of human data-entry errors, programming bugs, some transition effect or phase lag in stabilizing the data value to the true value.

The actual outliers need to be removed from the data set. For example, missing decimal in the cost of a toy US\$ 1.85 will make the cost 100 times more for a single toy. The result will thus be affected by a small or large amount. When valid data is identified as an outlier, then also the results are affected.

The statistical mean is computed from the product of each observed value v of *Values* with probability (or weight) P and then taking the average. The variance equals the difference of a value with respect to the mean, then square that, and then average the results. Standard deviation is just the square root of the variance.

Module 5

Chapter 9 (T2) – Text, Web Content and Link Analytics

• 9.1 Introduction

Text Analytics often termed as 'text mining' refers to analyzing and extracting the meanings, patterns, correlations and structure hidden in unstructured and semi-structured textual data. Text data stores consist of strong temporal dimensions, have modularity over time and sources, such as topics and sentiments.

Methods of machine-learning are prevalent in text analytics also. For example, when a user books an air-flight ticket using a tablet or desktop, the user receives an SMS on the mobile about details of the booking and flight timings. An ML algorithm, such as *Windows Crotona* at the mobile reads and learns by itself from the SMSs received at the phone. *Crotona* uses the ML for the SMS text analysis.. Learning results in SMS alerts to the user. An alert is reminder a day before the flight. Another alert is two hours before the flight, about the need to reach the airport. Those alerts are system-generated without prior request from the user.

The reader is required to know the meaning of the following select key terms:

Vector refers to an entity with number of interrelated elements. For example, a data point consists of n-elements in an n-dimensional space, and represents a vector to that point from the origin in the space. A word is a vector of characters as the elements. Consider a vector representation of word 'McGraw-Hill', then vector $V_{MH} = [M, c, G, r, a, w, -, H, i, l, l]$. V_{MH} is vector of 11 elements (characters) that refers to word McGraw-Hill.

Feature refers to a set of properties associated with an entity, object or category. For example, feature of properties, such as description of data analysis, data cleaning, data visualization and other topics in a book on data analytics.

Category refers to a classification on the basis of set of distinct features (for example, a category of text, document, cars, toys, students, news or fruits).

Dimension refers to a number of associated values, features or states, along the distinct spaces (dimensions). For example, a sentence has a number of words, each word has a number of characters, each word may have a feature, and so the sentences are in a three-dimensional space. Two dimensions are in metric spaces, which mean values in quantifiable spaces, such as the number of words, probability of occurrences in sentences, etc. Third dimension is in feature space, measured by a feature such as noun, verb, adverb, preposition, punctuation marks and stop word.

Another example is *apples*. Metric space variables are values, such as variables n , *number of apples* of specific properties, and P the *probability of preferring* apples of specific properties. Assume, feature space variables are four properties, *colour, shape, type* and *freshness*. Metric parameters and properties of apples are said to be in six-dimensional space, two are metric space (n and P) and four are feature space.

Graph data model refers to the data modelled by a set of entities. The entities identify by vertices V . A set of relations or associations identifies by edges E . An edge e represents a relation or association between two entities. Nodes represent the entities in the graph. The model also represents a hierarchy between the parent and children nodes.

Graph data network organization refers to a structure created by organizing entities or objects in a network, such as social network, business network and student network. A network organization means where persons or entities interconnect with each other, and have areas of common interest, business or study. A graph enables ease in traversing from one entity, person or web page link to another in the network by following a path. Web graph and social network graph enable such analysis. A graph network organization models the web and social networks. Examples of social networks are SlideShare, LinkedIn, Facebook and Twitter. The analytics of social networks finds the link ranks, clusters and correlations. The analytics discovers hubs and communities.

Web content mining refers to the discovery of useful information from web documents and services. Search engines use web content mining. A search provides the links of the required information to the user.

Hyperlinks refer to links mentioned in the contents that enable the retrieval of contents at web, file, object or resources repository.

Link analytics means web structure mining of hyperlinks between web documents. The analytics of links and analyzing them for metrics such as page ranks, clusters, correlations, hubs and communities.

Count triangles Algorithm is an algorithm that finds a number of triangular relationships among the nodes. Triangular relationships mean interrelations between each other.

Graph node centrality metric means the centrality of a node in reference to other nodes using certain metrics. Metrics used for centrality of a node are degree, closeness, betweenness or other characteristics of the node, such as rank, belief, potential, expectation, evidence, reputation or status.

Degree centrality of a node refers to the number of direct connections. Having more number of direct connections is not always a better metric. Better measure is the fact that the connection directs to significant results and tell how the nodes connect to the isolated node.

Betweenness centrality is a measure that provides the extent to which a node lies on paths between other nodes. A node with high betweenness signifies high influence over what flows in the network indicating importance of link and single point of failure.

Closeness centrality is the degree to which a node is near all other nodes in a network (directly or indirectly). It reflects the ability to access information through the network.

• 9.2 Text Mining

Today, large amounts of textual data is generated in computing applications. Text stream arriving continuously over time generates text data. For example, news articles, news reports, online comments on news, online traffic reports, corporate reports, web searches, and contents at social media discussion-forums (such as LinkedIn, Twitter and Facebook), short messages on phones, chat messages, transcripts of phone conversations, blogs and e-mails.

The abundance of textual data leads to problems which relate to their collection, exploration and ways of leveraging data. Textual data presents challenges for computing and storage requirements, consists of a strong temporal dimension, has modularity over time and have sources such as topics and sentiments. Examples of text processing techniques are clustering analysis, classifications, evolution analysis and event detection. Following subsections describe text mining in details:

9.2.1 Text Mining

Four definitions are:

1. "Text mining refers to the process of deriving high-quality information from text." (Wikipedia)
2. "Text mining is the process of discovering and extracting knowledge from unstructured data." (National Center of Text Mining –The University of Manchester¹)
3. "Text mining is the process of analyzing collections of textual contents in order to capture key concepts themes, uncover hidden relationships, and discover the trends without requiring that you know the precise words or terms that authors have used to express those concepts." (IBM²)
4. "Text mining is a technique which helps in revealing the patterns and relationships in large volumes of textual content that are not visible to the naked eye, leading to new business opportunities and improvements in processes." (Amazon BigData Official Blog³)

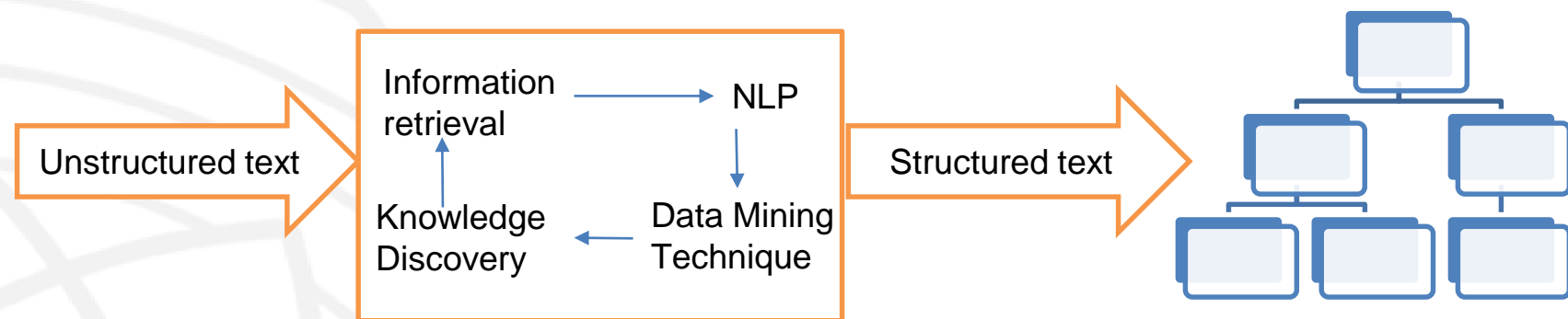
9.2.1.1 Text Mining Overview

Text mining includes extraction of high-quality information, discovering and extracting knowledge, and revealing patterns and relationships from unstructured data available in the form of text.

The term *text analytics* evolves from provisioning of strong integration with the already existing database technology, artificial intelligence, machine learning, data mining and text Data Store techniques. Information retrieval, natural language processing (NLP), classification, clustering and knowledge management are some of such useful techniques. Figure 9.1 shows process-pipeline in text-analytics.

9.2.1.2 Areas and Applications of Text Mining

Natural Language Processing (NLP) is a technique for analyzing, understanding and deriving meaning from human language. NLP involves the computer's understanding and manipulation of human language. NLP algorithms are typically based on ML algorithms. They automatically learn the rules. First, they analyze set of examples from a large collection of sentences in a book. Then, they make the statistical inferences.



NLP contributes to the field of human computer interaction by enabling several real-world applications such as automatic text summarization, sentiment analysis, topic extraction, named entity recognition, parts-of-speech tagging, relationship extraction and stemming. The common uses of NLP include text mining, machine translation and automated question answering.

Information Retrieval (IR) is a process of searching and retrieving a subset of documents from the abundant collection of documents. IR can also be defined as extraction of information required by a user. IR is an area derived fundamentally from database technology. One of the most popular applications of IR is searching the information on the web. Search engines provide IR using various advance techniques. For example, the crawler program is capable of retrieving information from a wide variety of data sources. Search methods use metadata or full-text indexing.

Information Extraction (IE) is a process in which the software extracts structured information from unstructured and/or semi-structured documents. IE finds the relationship within text or desired contents from text. IE ideally derives from machine learning, more specifically from the NLP domain. Content extraction from the images, audio or video is an example of information extraction.

IE requires a dictionary of extraction patterns (For example, “Citizen of <x>”, or “Located in <x>”) and a semantic lexicon (dictionary of words with semantic category labels).

Document Clustering is an application which groups text documents into clusters. Automating document organization, topic extraction and fast information retrieval or filtering use the document clustering method. For example, web document clustering facilitates easy search by users.

Document Classification is an application to classify text documents into classes or categories. The application is useful for publishers, news sites, blogs or areas where lot of contents are present.

Web Mining is an application of data mining techniques. They discover patterns from the web Data Store. The patterns facilitate understanding. They improve the services of web-based applications. Data mining of web usage provides the browsing behavior of a website.

Concept Extraction is an application that deals with the extraction of concept from textual data. Concept extraction is an area of text classification in which words and phrases are classified into a semantically similar group.

9.2.1.3 Text Mining Process

Text is most commonly used for information exchange. Unlike data stored in databases, text is unstructured, ambiguous and difficult to process. Text mining is the process that analyzes a text to extract information useful for a specific purpose.

Syntactically, a text document comprises characters that form words, which can be further combined to generate phrases or sentences. Text mining steps are (i) recognizing, extracting and using the information present in words. Along with searching of words, mining involves search for semantic patterns as well.

Text mining process consists of a process-pipeline. The pipeline processes execute in several phases. Mining uses the iterative and interactive processes. The processing in pipeline does text mining efficiently and mines the new information. Figure 9.2 shows five phases of the process pipeline.

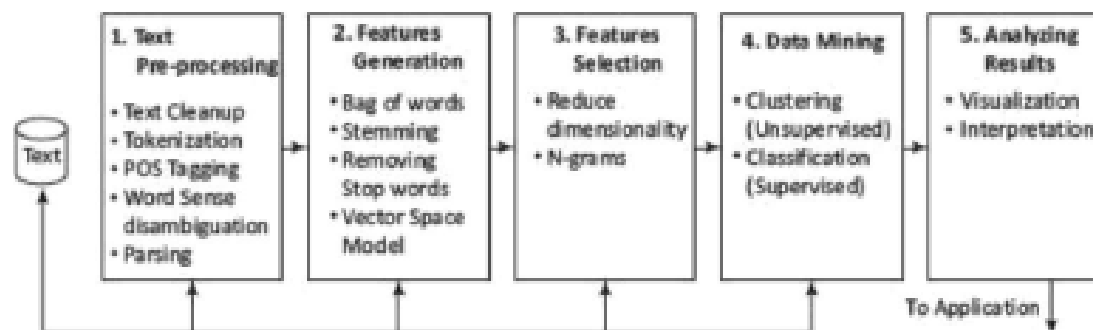


Figure 9.2 Five phases in a process pipeline

The five phases for processing text are as follows:

Phase 1: Text pre-processing enables Syntactic/Semantic text-analysis and does the followings:

1. *Text cleanup* is a process of removing unnecessary or unwanted information. Text cleanup converts the raw data by filling up the missing values, identifies and removes outliers, and resolves the inconsistencies. For example, removing comments, removing or escaping “%20” from URL for the web pages or cleanup the typing error, such as teh (the), do n’t (do not) [%20 specifies space in a URL].
2. *Tokenization* is a process of splitting the cleanup text into tokens (words) using white spaces and punctuation marks as delimiters.
3. *Part of Speech (POS) tagging* is a method that attempts labeling of each token (word) with an appropriate POS. Tagging helps in recognizing names of people, places, organizations and titles. English language set includes the noun, verb, adverb, adjective, prepositions and conjunctions. Part of Speech encoded in the annotation system of the Penn Treebank Project has 36 POS tags.⁴
4. *Word sense disambiguation* is a method, which identifies the sense of a word used in a sentence; that gives meaning in case the word has multiple meanings. The methods, which resolve the ambiguity of words can be context or proximity based. Some examples of such words are bear, bank, cell and bass.
5. *Parsing* is a method, which generates a parse-tree for each sentence. Parsing attempts and infers the precise grammatical relationships between different words in a given sentence.

Phase 2: Features Generation is a process which first defines features (variables, predictors). Some of the ways of feature generations are:

1. *Bag of words*—Order of words is not that important for certain applications.

Text document is represented by the words it contains (and their occurrences). Document classification methods commonly use the bag-of-words model. The pre-processing of a document first provides a document with a bag of words. Document classification methods then use the occurrence (frequency) of each word as a feature for training a classifier. Algorithms do not directly apply on the bag of words, but use the frequencies.

2. *Stemming*—identifies a word by its root.

- (i) Normalizes or unifies variations of the same concept, such as *speak* for three variations, i.e., speaking, speaks, speakers denoted by [speaking, speaks, speaker \rightarrow speak]
- (ii) Removes plurals, normalizes verb tenses and remove affixes.

Stemming reduces the word to its most basic element. For example, impurification \rightarrow pure.

3. *Removing stop words* from the feature space—they are the common words, unlikely to help text mining. The search program tries to ignore stop words. For example, ignores *a, at, for, it, in* and *are*.
4. *Vector Space Model (VSM)*—is an algebraic model for representing text documents as vector of identifiers, word frequencies or terms in the document index. VSM uses the method of term frequency-inverse document frequency (TF-IDF) and evaluates how important is a word in a document.

When used in document classification, VSM also refers to the bag-of-words model. This bag of words is required to be converted into a term-vector in VSM. The term vector provides the numeric values corresponding to each term appearing in a document. The term vector is very helpful in feature generation and selection.

Phase 3: Features Selection is the process that selects a subset of features by rejecting irrelevant and/or redundant features (variables, predictors or dimension) according to defined criteria. Feature selection process does the following:

1. *Dimensionality reduction*—Feature selection is one of the methods of division and therefore, dimension reduction. The basic objective is to eliminate irrelevant and redundant data. Redundant features are those, which provide no extra information. Irrelevant features provide no useful or relevant information in any context.

Principal Component Analysis (PCA) and Linear Discriminate Analysis (LDA) are dimension reduction methods. Discrimination ability of a feature measures relevancy of features. Correlation helps in finding the redundancy of the feature. Two features are redundant to each other if their values correlate with each other.

2. *N-gram evaluation*—finding the number of consecutive words of interest and extract them. For example, 2-gram is a two words sequence, ["tasty food", "Good one"]. 3-gram is a three words sequence, ["Crime Investigation Department"].
3. *Noise detection and evaluation of outliers* methods do the identification of unusual or suspicious items, events or observations from the data set. This step helps in cleaning the data.

The feature selection algorithm reduces dimensionality that not only improves the performance of learning algorithm but also reduces the storage requirement for a dataset. The process enhances data understanding and its visualization.

Phase 4: Data mining techniques enable insights about the structured database that resulted from the previous phases. Examples of techniques are:

1. Unsupervised learning (for example, clustering)
 - (i) The class labels (categories) of training data are unknown
 - (ii) Establish the existence of groups or clusters in the data

Good clustering methods use high intra-cluster similarity and low inter-cluster similarity. Examples of uses – blogs, patterns and trends.

2. *Supervised learning (for example, classification)*

- (i) The training data is labeled indicating the class
- (ii) New data is classified based on the training set

Classification is correct when the known label of test sample is identical with the resulting class computed from the classification model.

Examples of uses are *news filtering application*, where it is required to automatically assign incoming documents to pre-defined categories; *email spam filtering*, where it is identified whether incoming email messages are spam or not.

Example of text classification methods are *Naïve Bayes Classifier* and *SVMs*.

3. *Identifying evolutionary patterns in temporal text streams*—the method is useful in a wide range of applications, such as summarizing of events in news articles and extracting the research trends in the scientific literature.

Phase 5: Analysing results

- (i) Evaluate the outcome of the complete process.
- (ii) Interpretation of Result– If acceptable then results obtained can be used as an input for next set of sequences. Else, the result can be discarded, and try to understand what and why the process failed.
- (iii) Visualization – Prepare visuals from data, and build a prototype.
- (iv) Use the results for further improvement in activities at the enterprise, industry or institution.

9.2.1.5 Text Mining Challenges

The challenges in the area of text mining can be classified on the basis of documents area-characteristics. Some of the classifications are as follows:

1. NLP issues:
 - (i) POS Tagging
 - (ii) Ambiguity
 - (iii) Tokenization
 - (iv) Parsing
 - (v) Stemming
 - (vi) Synonymy and polysemy
2. Mining techniques:
 - (i) Identification of the suitable algorithm(s)
 - (ii) Massive amount of data and annotated corpora
 - (iii) Concepts and semantic relations extraction
 - (iv) When no training data is available
3. Variety of data:
 - (i) Different data sources require different approaches and different areas of expertise
 - (ii) Unstructured and language independency
4. Information visualization
5. Efficiency when processing real-time text stream
6. Scalability

9.2.1.6 Supervised text Classification

The supervised text classification requires labeled documents and additional knowledge from experts. The algorithms exploit the training data (where zero or more categories) to learn a classifier, which classifies new text documents and labels each document. A document is considered as a positive example for all categories with which it is labeled, and as a negative example to all others. The task of a training algorithm for a text classifier is to find a weight vector which best classifies new text documents.

The different approaches for supervised text classification are:

- (i) K-Nearest Neighbour Method
- (ii) Support Vector Machine
- (iii) Naïve Bayes Method
- (iv) Decision Tree
- (v) Decision Rule

K-Nearest Neighbours (KNN) method makes use of training text document. The training documents are the previously categorized set of documents. They train the system to understand each category. The classifier uses the training 'model' to classify new incoming documents. KNN assumes that close-by objects are more probable in the same category. KNN finds k objects in the large number of text documents, which have most similar query responses. Thus, in KNN, predictions are based on a method that is used to predict new (not observed earlier) text data. The predictions are by (i) majority vote method (for classification tasks) and (ii) averaging (for regression) method over a set of K-nearest examples.

The decision trees or decision rules are built to predict the category for an input document. A decision tree or rule represents a set of nested logical if-then conditions on the observed values of the text features that enable the prediction of the target variable. The decision tree and decision rules are also used to classify (categorize) the document. Classification is done by recursively splitting the text features into a set of non-overlapping regions (Refer Section 6.8). (Section 6.7.4)

9.3 ↓ WEB MINING, WEB CONTENT AND WEB USAGE ANALYTICS

Web is a collection of interrelated files at web servers. Web data refers to

(i) web content—text, image and records, (ii) web structure—hyperlinks and tags, and (iii) web usage—http logs and application server logs.

Features of web data are:

1. Volume of information and its ready availability
2. Heterogeneity
3. Variety and diversity (Information on almost every topic is available using different forms, such as text, structured tables and lists, images, audio and video.)
4. Mostly semi-structured due to the nested structure of HTML code
5. Hyperlinks among pages within a website, and across different websites
6. Redundant or similar information may be present in several pages
7. Mostly, the web page has multiple sections (divisions), such as main contents of the page, advertisements, navigation panels, common menu for all the pages of a website and copyright notices
8. A web form or HTML form on a web page enables a user to enter data that is sent to a server for processing
9. Website contents are dynamic in nature where information on the web pages constantly changes, and fast information growth takes place such as conversations between users, social media, etc.

LO 9.2

Mining the web-links, web-structure and web-contents, and analyzing the web graphs

9.3.1 Web Mining

Data Mining is a process of discovering patterns in large datasets to gain knowledge. The process can be shown as [Raw Data → Patterns → Knowledge]. Web data mining is the mining of web data. Web mining methods are in multidisciplinary domains: (i) data mining, ML, natural language, (ii) processing, statistics, databases, information retrieval, and (iii) multimedia and visualization.

Web consists of rich features and patterns. A challenging task is retrieving interesting content and discovering knowledge from web data. Web offers several opportunities and challenges to data mining.

Definition of Web Mining

Web mining refers to the use of techniques and algorithms that extract knowledge from the web data available in the form of web documents and services. Web mining applications are as follows:

- (i) Extracting the fragment from a web document that represents the full web document
- (ii) Identifying interesting graph patterns or pre-processing the whole web graph to come up with metrics, such as PageRank
- (iii) User identification, session creation, malicious activity detection and filtering, and extracting usage path patterns

Web Mining Taxonomy

Web mining can broadly be classified into three categories, based on the types of web data to be mined. Three ways are web content mining, web structure mining and web usage mining. Figure 9.6 shows the taxonomy of web mining.

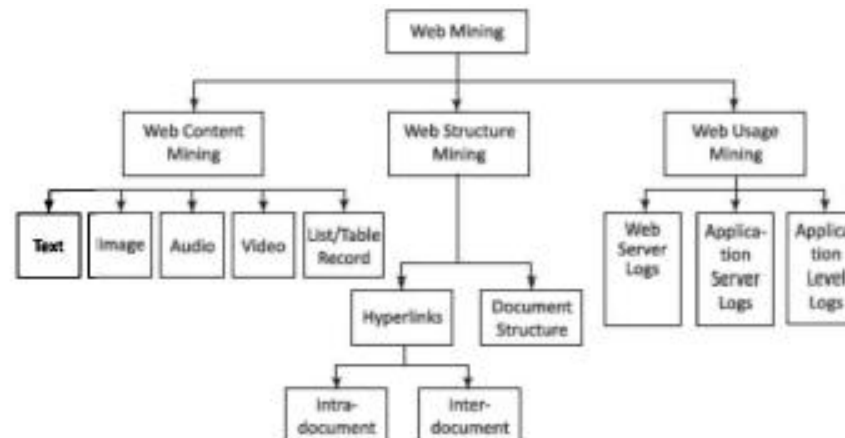


Figure 9.6 Web mining taxonomy

Web content mining is the process of extracting useful information from the contents of web documents. The content may consist of text, images, audio, video or structured records, such as lists and tables.

Web structure mining is the process of discovering structure information from the web. Based on the kind of structure-information present in the web resources, web structure mining can be divided into:

1. **Hyperlinks:** the structure that connects a location at a web page to a different location, either within the same web page (intra-document hyperlink) or on a different web page (inter-document hyperlink)
2. **Document Structure:** The structure of a typical web graph consists of web pages as nodes, and hyperlinks as edges connecting the related pages.

Web usage mining is the application of data mining techniques which discover interesting usage patterns from web usage data. The data contains the identity or origin of web users along with their browsing behavior at a web site. Web usage mining can be classified as:

- (i) **Web Server logs:** Collected by the web server and typically include IP address, page reference and access time.
- (ii) **Application Server Logs:** Application servers typically maintain their own logging and these logs can be helpful in troubleshooting problems with services.
- (iii) **Application Level Logs:** Recording events usually by application software in a certain scope in order to provide an audit trail that can be used to understand the activity of the system and to diagnose problems.

9.3.2 Web Content Mining

Web Content Mining is the process of information or resource discovery from the content of web documents across the World Wide Web. Web content mining can be (i) direct mining of the contents of documents or (ii) mining through search engines. They search fast compared to direct method.

Web content mining relates to both, data mining as well as text mining. Following are the reasons:

- (i) The content from web is similar to the contents obtained from database, file system or through any other mean. Thus, available data mining techniques can be applied to the web.
- (ii) Content mining relates to text mining because much of the web content comprises texts.
- (iii) Web data are mainly semi-structured and/or unstructured, while data mining is structured and the text is unstructured.

Applications

Following are the applications of content mining from web documents:

1. Classifying the web documents into categories
2. Identifying topics of web documents
3. Finding similar web pages across the different web servers
4. Applications related to relevance:
 - (a) Recommendations – List of top “n” relevant documents in a collection or portion of a collection
 - (b) Filters – Show/Hide documents based on some criterion
 - (c) Queries – Enhance standard query relevance with user, role, and/or task-based relevance

9.3.2.1 Common Web Content Mining Techniques

Pre-processing of contents The pre-processing steps are quite similar to the pre-processing for text mining. The content preparation involves:

1. **Extraction of text from HTML**
2. **Data cleaning** by filling up the missing values and smoothing the noisy data
3. **Tokenizing**: Generates the tokens of words from the cleaned up text
4. **Stemming**: Reduce the words to their roots. The different grammatical forms or declinations of verbs identify and index (count) as the same word. For example, stemming will ensure that both “closed” and “closing” are derived from the same word “close”. Stemming algorithm, *Porter*, can be used here. The java code for *Porter* stemming algorithm can be obtained from <https://tartarus.org/martin/PorterStemmer/java.txt>,
5. **Removing the stop words**: The common words unlikely to help in the mining process such as articles (a, an, the), or prepositions (such as, to, in, for) are removed.
6. **Calculate collection wide-word frequencies**: The distinct-word stem obtained after stemming process and removing the stop words results into a list of significant words (or terms). Calculating the occurrence of a significant term (t) in a collection is called collection frequency (CF_t). CF counts the multiple occurrences.)

Now, find the number of documents in the collection that contains the specific term (t). This numeric measure is the document frequency (DF_t).

7. **Calculate per Document Term Frequencies (TF)**. TF is a numeric measure that is used to score the importance of a word in a document based on how often it appeared in that document (Refer Example 9.1).
8. **Bag of words**: Web document is represented by the words it contains (and their occurrences).

The following example explains the concept of CF and DF using the data of toy sales collection.

Mining Tasks for Web Content Analytics

Following are the tasks for web content analytics:

1. **Classification** – A supervised technique which:

- (i) Identifies the class or category a new web documents belongs to from the set of predefined classes or categories
- (ii) Categories in the form of a term vector that are produced during a “training” phase
- (iii) Employs algorithms using term vector to categorize the new data according to the observations at the training set.

2. **Clustering** – An unsupervised technique:

- (i) Groups the web documents (clustered) with similar features using some similarity measure
- (ii) Uses no pre-defined perception of what the groups should be
- (iii) Measures most common similarity using the dot product between two web document vectors.

3. **Identifying the association** between web documents – Association rules help to identify correlation between web pages that occur mostly together.

The other significant mining tasks are:

1. **Topic identification, tracking and drift analysis** – A way of organizing the large amount of information retrieved from the web is categorizing the web pages into distinct topics. The categorization can be based on a similarity metric, which includes textual information and co-citation relations. Clustering or classification techniques can automatically and effectively identify relevant topics and add them in a topic-wise collection library.

Adding a new document to a collection library includes:

- (i) Assigning each document to an existing topic (category)
- (ii) Re-checking of collection for the emergence of new topics
- (iii) Tracking the number of views to a collection
- (iv) Identifying the drift in a topic(s)

2. *Concept hierarchy creation* - Concept hierarchy is an important tool for capturing the general relationship among web documents. Creation of concept hierarchies is important to understand a category and sub-categories to which a document belongs. The clustering algorithms leverage more than two clusters, which merge into a cluster. That is merging the sub-clusters into a cluster.

Important factors for creation of concept hierarchy include:

- (i) Identifying the organization of categories, such as flat, tree or network
- (ii) Planning the maximum number of categories per document
- (iii) Building category dimensions, such as domain, location, time, application and privileges.

3. *Relevance of content* - Relevance or the applicability of web content can be measured with respect to any of the following basis:

- (i) Document relevance describes the usefulness of a given document in a specified situation.
- (ii) Query-based relevance is the most useful method to assess the relevance of web pages. Query-based relevance is used in information retrieval tools such as search engines. The method calculates the similarity between query (search) keywords and document. Similarity, results can be refined through additional information such as popularity metric as seen in Google or the term positions in AltaVista.
- (iii) User-based relevance is useful in personal aspects. User profiles are maintained, and similarity between the user profile and document is calculated. The relevance is often used in push notification services.
- (iv) Role/task-based relevance is quite similar to user-based relevance. Instead of a user, here the profile is based on a particular role or task. Multiple users can provide input to profile.

9.3.3 Web Usage Mining

Web usage mining discovers and analyses the patterns in click streams. Web usage mining also includes associated data generated and collected as a consequence of user interactions with web resources.

Figure 9.7 shows three phases for web usage mining.

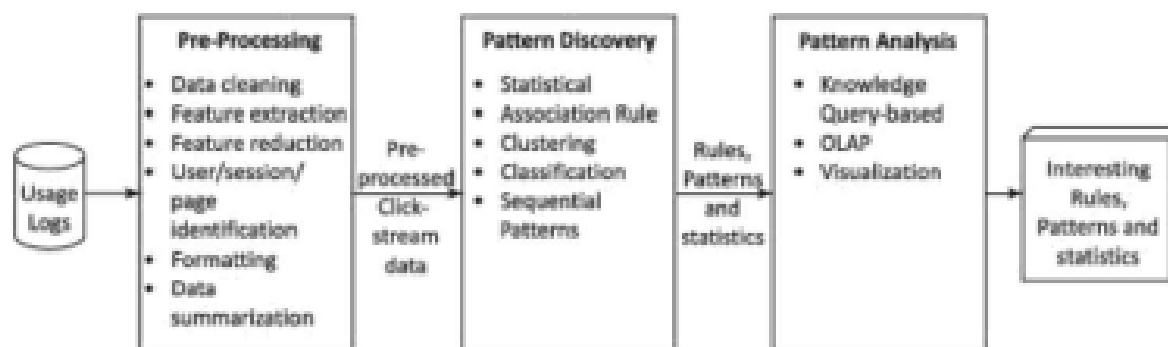


Figure 9.7 Process of web usage mining

The phases are:

1. Pre-processing – Converts the usage information collected from the various data sources into the data abstractions necessary for pattern discovery.
2. Pattern discovery – Exploits methods and algorithms developed from fields, such as statistics, data mining, ML and pattern recognition.
3. Pattern analysis – Filter out uninteresting rules or patterns from the set found during the pattern discovery phase.

Usage data are collected at server, client and proxy levels. The usage data collected at the different sources represent the navigation patterns of the overall web traffic. This includes single-user, multi-user, single-site access and multi-site access patterns.

9.3.3.1 Pre-processing

The common data mining techniques apply on the results of pre-processing using vector space model (Refer Example 9.2).

Pre-processing is the data preparation task, which is required to identify:

- (i) User through cookies, logins or URL information
- (ii) Session of a single user using all the web pages of an application
- (iii) Content from server logs to obtain state variables for each active session
- (iv) Page references.

The subsequent phases of web usage mining are closely related to the smooth execution of data preparation task in pre-processing phase. The process deals with (i) extracting of the data, (ii) finding the accuracy of data, (iii) putting the data together from different sources, (iv) transforming the data into the required format and (iv) structure the data as per the input requirements of pattern discovery algorithm.

Pre-processing involves several steps, such as data cleaning, feature extraction, feature reduction, user identification, session identification, page identification, formatting and finally data summarization.

9.3.3.2 Pattern Discovery

The pre-processed data enable the application of knowledge extraction algorithms based on statistics, ML and data mining algorithms. Mining algorithms, such as path analysis, association rules, sequential patterns, clustering and classification enable effective processing of web usages. The choice of mining techniques depends on the requirement of the analyst. Pre-processed data of the web access logs transform into knowledge to uncover the potential patterns and are further provided to pattern analysis phase.

Some of the techniques used for pattern discovery of web usage mining are:

Statistical techniques They are the most common methods which extract the knowledge about users. They perform different kinds of descriptive statistical analysis (frequency, mean, median) on variables such as page views, viewing time and length of path for navigational.

Statistical techniques enable discovering:

- (i) The most frequently accessed pages
- (ii) Average view time of a page or average length of a path through a site
- (iii) Providing support for marketing decisions

Association rule The rules enable relating the pages, which are most often referenced together in a single server session. These pages may not be directly connected to one another using the hyperlinks.

Other uses of association rule mining are:

- (i) Reveal a correlation between users who visited a page containing similar information. For example, a user visited a web page related to admission in an undergraduate course to those who search an eBook related to any subject.
- (ii) Provide recommendations to purchase other products. For example, recommend to user who visited a web page related to a book on data analytics, the books on ML and Big Data analytics also.
- (iii) Provide help to web designers to restructure their websites.
- (iv) Retrieve the documents in prior in order to reduce the access time when loading a page from a remote site.

Clustering is the technique that groups together a set of items having similar features. Clustering can be used to:

- (i) Establish groups of users showing similar browsing behaviors
- (ii) Acquire customer sub-groups in e-commerce applications
- (iii) Provide personalized web content to users
- (iv) Discover groups of pages having related content. This information is valuable for search engines and web assistance providers.

Classification The method classifies data items into predefined classes. Classification is useful for:

- (i) Developing a profile of users belonging to a particular class or category
- (ii) Discovery of interesting rules from server logs. For example, 3750 users watched a certain movie, out of which 2000 are between age 18 to 23 and 1500 out of these lives in metro cities.

Classification can be done by using supervised inductive learning algorithms, such as decision tree classifiers, Naïve Bayesian classifiers, k-nearest neighbour classifiers, support vector machines.

Sequential pattern discovery User navigation patterns in web usage data gather web page trails that are often visited by users in the order in which pages are visited. Markov Model can be used to model navigational activities in the website. Every page view in this model can be represented as a state. Transition probability between two states can represent the probability that a user will navigate from one state to the other. This representation allows for the computation of a number of significant user or site metrics that can lead to useful rules, pattern, or statistics.

9.3.3.3 Pattern Analysis

The objective of pattern analysis is to filter out uninteresting rules or patterns from the rules, patterns or statistics obtained in the pattern discovery phase.

The most common form of pattern analysis consists of:

- (i) A knowledge query mechanism such as SQL
- (ii) Another method is to load usage data into a data cube in order to perform Online Analytical Processing (OLAP) operations
- (iii) Visualization techniques, such as graphing patterns or assigning the colors to different values, can often highlight overall patterns or trends in the data
- (iv) Content and structure information can filter out patterns containing pages of a certain usage type, content type or pages that match a certain hyperlink structure.

Data cube enables visualizing data from different angles. For example, toys data visualization using category, colour and children preferences. Another example, news from category, such as sports, success stories, films or targeted readers (children, college students, etc).

9.4 | PAGE RANK, STRUCTURE OF WEB AND ANALYZING A WEB GRAPH

LO 9.3

PageRanking, analysis of web-structure, and discovering hubs, authorities and communities in web-structure

Sections 9.2 and 9.3 described text data and web contents analysis. Hyperlinks links exist between the web contents. Link analysis finds the answers to the following:

1. Can a linked (web) page rank them higher or lower?
2. Can the links be modeled as edges of graphs, structure of web as graph network, and applied the tools *same as for graph analytics*?
3. Can web graph mining method analyze and find a link sending spams?
4. Does a set of links correspond to a hub? Do the links correspond to an authority?
5. Does a linked page has higher authority compared to others?

Links analysis applies to domains of social networks and e-mail. The following sub-sections describe the applications of link analysis:

9.4.1 Page Rank Definition

The in-degree (visibility) of a link is the measure of number of in-links from other links. The out-degree (luminosity) of a link is number of other links to which that link points.

PageRank definition according to earlier approaches

Assume a web structure of hyperlinks. Each hyperlink in-links to a number of hyperlinks and out-links to a number of pages. A page commanding higher authority (rank) has greater number of in-degrees than out-degrees. Therefore, one measure of a page authority can be in-degrees with respect to out-degrees.

PageRank refers to the authority of the page measured in terms of number of times a link is sought after.

PageRank definition according to the new approach

Earlier approach of page ranking based on in-links and out-links does not capture the relative authority (importance) of the parents. Page and co-authors (1998) defined a page ranking method,⁵ which considers the entire web in place of local neighbourhood of the pages and considers the relative authority of the parent links (over children).

9.4.2 Web Structure

Web structure models as directed-graphs network-organization. Vertex of the directed graph models an anchor. Let n = number of hyperlinks at the page U . Assume \mathbf{u} is a vector with elements u_1, u_2, \dots, u_n . Each page $Pg(\mathbf{u})$ has anchors, called hyperlinks. Page $Pg(\mathbf{v})$ consists of text document with m number of hyperlinks. \mathbf{v} is a vector with elements v_1, v_2, \dots, v_m . The m is number of hyperlinks at $Pg(\mathbf{v})$. A vertex u directs to another Page V . A page $Pg(\mathbf{v})$ may have number of hyperlinks directed by out-edges to other page $Pg(\mathbf{w})$. Consider the following hypotheses:

1. Text at the hyperlink represents the property of a vertex u that describes the destination V of the out-going edge.
2. A hyperlink in-between the pages represents the conferring of the authority.

Pages U and U' hyperlinks u and u' out-linking to Page V . Let Page U has three hyperlinks parenting three Pages, V one, W two, X two, U' one, and Y two, respectively. Figure 9.8 shows a web structure consisting of pages and hyperlinks.

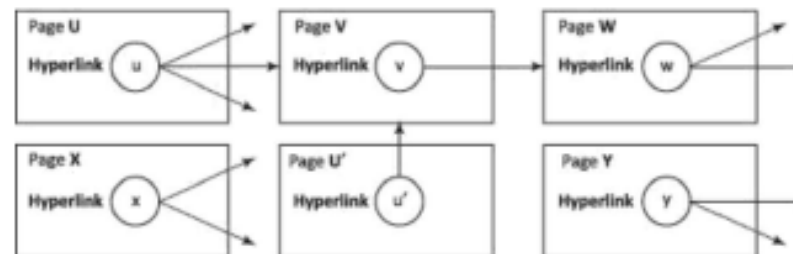


Figure 9.8 Web structure with hyperlinks from a parent to one or more pages

9.4.2.1 Dead Ends

Dead-end web pages refer to pages with no out-links. When a web page links to such pages, its page rank gets reduced. Dead ends are on a website having poor linking structure.

The web structure of service pages may have pages with a dead end. The end causes no further flows for further action and no internal links. Good website structures have the pages designed such that they specifically gently guide the visitors toward actions and towards next step. For example, if one searches for a book title on Amazon, then visitor gets links of other books also on a similar topic.

9.4.2.2 Analyzing and Implementing a System with Web Graph Mining

Number of metrics analyze a system using web graph mining. Following are the examples:

1. In-degrees and out-degrees
2. Closeness is centrality metric. Closeness, $Cc(v) = 1 / \sum_{u \neq v} gdist(v, u)$, where $gdist$ is the geodesic distance between vertex v with u and sum is over all u linked with V . Geodesic distance means the number of edges in a shortest path connecting two vertices. Assume v has an edge with w , and w has an edge with u . Assume u does not have direct edge from v . Then, geodesic distance = 2 (two edges between v and u in shortest path).
3. Betweenness
4. PageRank and LineRank
5. Hubs and authorities
6. Communities parameters, triangle count, clustering coefficient, K-neighbourhood
7. Top K-shortest paths

9.4.3 Computation of PageRank and PageRank Iteration

Assume that a web graph models the web pages. Page hyperlinks are the property of the graph node (vertex). Assume a Page, $Pg(v)$ in-links from $Pg(u)$, and $Pg(u)$ out-linking similar to $Pg(v)$, to total $N_{out}[(Pg(u))]$ pages. Figure 9.9 shows $Pg(v)$ in-links from $Pg(u)$ and other pages.

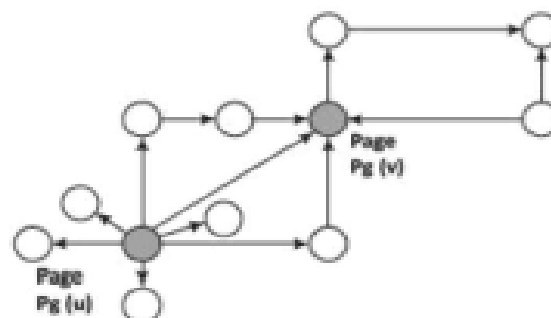


Figure 9.9 Page $Pg(v)$ in-links from $Pg(u)$ and other pages

N_{out} for page U is 7 and for V is 1 in the figure. Number of in-linking N_{in} for page V is 4. Two algorithms to compute page rank are as follows:

1. PageRank algorithm using the in-degrees as conferring authority

Assume that the page U, when out-linking to Page V “considers” an equal fraction of its authority to all the pages it points to, such as Pgv. The following equation gives the initially suggested page rank, PR (based on in-degrees) of a page Pgv:

$$PR(P_{gv}) = nc \cdot \sum_{P_{gu}: P_{gu} \rightarrow P_{gv}} [PR(P_{gu})/N(P_{gu})] \quad (9.21)$$

where N(Pgu) is the total number of out-links from U. Sum is over all Pgv in-links. Normalization constant denotes by nc, such that PR of all pages sums equal to 1.

However, just measuring the in-degree does not account for the authority of the source of a link. Rank is flowing among the multiple sets of the links. When Pgv in-links to a page Pgu, its rank increases and when page Pgu out-links to other new links, it means that N (Pgu) increases, then rank PR(Pgv) sinks (decreases). Eventually, the PR (Pgv) converges to a value.

Therefore, rank computation algorithm iterates the rank flowing computations as shown below:

EXAMPLE 9.7

Assume S corresponds to a set of pages. Initialize $\forall P_g \in S$. Symbols mean that initialize all pages Pg contained in the S and initialize Page Rank (Pgv) for each page as follows:

$$PR_{init}(P_{gv}) = 1/|S| \quad (9.22)$$

How are the page ranks of the pages in a given set of pages iterated and computed till the ranks do not change (within specified margin, that means until converge)?

SOLUTION

Iterate and compute PR (Pgv) for each page as follows:

Until ranks do not change (within specified margin) (that means converge)

{

for each Pgv $\in S$ compute,

$$PR(P_{gv}) = \sum_{P_{gu}: P_{gu} \rightarrow P_{gv}} [PR(P_{gu}) / N(P_{gu})] \quad (9.22)$$

and normalization constant,

$$nc = \sum_{P_{gv} \in S} [PR_{new}(P_{gv})] \quad (9.23a)$$

$$\text{for each } P_{gv} \in S: PR(P_{gv}) = nc \cdot PR(P_{gv}) \quad (9.23b)$$

}

2. PageRank algorithm using the relative authority of the parents over linked children

A method of PageRank considers the entire web in place of local neighbourhood of the pages and considers the relative authority of the parents (children). The algorithm uses the relative authority of the parents (children) and adds a rank for each page from a rank source.

The PageRank method considers assigning weight according to the rank of the parents. Page rank is proportional to the weight of the parent and inversely proportional to the out-links of the parent.

Assume that (i) Page v (P_{gv}) has in-links with parent Page u (P_{gu}) and other pages in set $PA(v)$ of parent pages to v that means $\in PA(v)$, (ii) $R(v)$ is PageRank of P_{gv} , (iii) $R(u)$ is weight (importance/rank) of P_{gu} , and (iv) $ch(u)$ is weight of child (out-links) of P_{gu} . Then the following equation gives PageRank $R(v)$ of link v :

$$R(v) = \sum_{u \in PA(v)} \left[R(u) / [ch(u)] \right] \quad (9.25)$$

where $PA(v)$ is a set of links who are parents (in-links) of link v . Sum is over all parents of v . nc is normalization constant whose sum of weights is 1.

Assume that a rank source E exists that is addition to the rank of each page $R(v)$ by a fixed rank value $E(v)$ for P_{gv} . $E(v)$ is fraction α of $[1/|PA(v)|]$.

An alternative equation is as follows:

$$R(v) = nc \cdot \left\{ (1 - \alpha) \sum_{u \in PA(v)} \left[\frac{R(u)}{[ch(u)]} \right] + \alpha \cdot E(v) \right\}. \quad (9.26)$$

where $nc = [1/R(v)]$. $R(v)$ is iterated and computed for each parent in the set $PA(v)$ till new value of $R(v)$ does not change within the defined margin, say 0.001 in the succeeding iterations.

Significance of a PageRank can be seen as modeling a “random surfer” that starts on a random page and then at each point: $E(v)$ models the probability that a random link jumps (surfs) and connect with out-link to Pgv . $R(v)$ models the probability that the random link connects (surf) to Pgv at any given time. The addition of $E(v)$ solves the problem of Pgv by chance out-linking to a link with dead end (no outgoing links).

Therefore, rank computation algorithm iterates the rank flowing computations as shown in Example 9.8.

EXAMPLE 9.8

Assume PA corresponds to a set of parent pages to a page v . Initialize $\forall Pg \in PA(v)$. Symbols mean that initialize all pages u contained in the set of parent pages of $PA(v)$ and initialize Page Rank $R(v)$ for each page as follows:

$$R(v) = [1/|PA(v)|]$$

How are the page ranks of the pages in a given set of pages iterated and computed till the ranks do not change (within specified margin, that means untill converges)?

SOLUTION

Iterate and compute $R(v)$ for each page as follows:

Until ranks do not change that means converges (within specified margin, say 0.001)

for each $v \in PA(v)$ compute,

$$R(v) = nc \cdot \left\{ (1-\alpha) \sum_{u \in PA(v)} \left[\frac{R(u)}{|ch(u)|} \right] + \alpha \cdot E(v) \right\} \quad (9.27)$$

and normalization constant,

$$nc = \sum_{u \in PA(v)} [R(v)] \quad (9.28a)$$

$$\text{for each } v \in PA(v): R(v) = nc \cdot R(v) \quad (9.28b)$$

}

PageRank Iteration using MapReduce functions in Spark Graph

The computation of PageRank using SparkGraph method (Section 8.5),

```
graph.pageRank(0.0001).vertices  
  
ranksByUsername = users.join(ranks).map{case id, (username, rank)} => (username, rank).
```

The method includes conversions to MapReduce functions and using HDFS compatible files. Functions PageRank (), ranksByUsername () do the computations using the PageRankObject. GraphX consists of these functions (GraphOps). GraphX Operators includes the functions (Section 8.5).

Static PageRank algorithm runs for a fixed number of iterations, while dynamic PageRank runs until the computed rank converges. Convergence means that after certain iterations, the rank does not change significantly and any change remains within a pre-specified tolerance. Thereafter the iterations stop.

Assume specified tolerance at the start of iterations is 0.0001 (1 in 10000). When the rank does not change beyond that tolerance, it means rank value will converge and then the iterative process will stop.

9.4.4 Topic Sensitive PageRank and Link Spam

Number of methods have been suggested for computations of topic-sensitive page ranking, R_{TS} . The $R_{TS}(v)$ of a page $P(v)$ may be higher for a specific topic compared to other topics. A topic associates with a distinct bag of words for which the page has higher probability of surfing than other bags for that topic.

Topic-sensitive PageRank method uses surfing weights (probabilities) for the pages containing the topic or bag of words corresponding to a topic. Method for creating topic-sensitive PageRank is to compute the bias to rank $R(v)$ and thus increase the effect of certain pages containing that topic or bag of words.

Refer equation (9.25) for computations of $R(v)$, and equation (9.26) for computations after introducing additional influence to the page. A method of introducing biasing is simple. It assumes that a rank source E exists that is additional having in-links from other pages, and thus adds to the rank of each page $R(v)$ by a fixed (uniform) or non-uniform weight factor α . The factor α is a multiplication factor to actual in-links without the bias.

Recapitulate equation (9.26). Probability of random jump to page v is $E(v)$. An alternative equation for topic sensitive PageRank, $R(v)$ computation for page $P(v)$ is as follows:

$$R(v) = nc \cdot \left\{ (1 - \alpha_t) \cdot P(v) \sum_{u \in \text{In}(v)} \left[\frac{R(u)}{|\text{ch}(u)|} \right] + \alpha_t \cdot E(v) \right\}. \quad (9.29)$$

Probability of random jump to page v is $E(v)$. $\alpha_t = 0$ for page unrelated to a topic α is not 0 for page related to a topic. α_t = surfing probability for in-links for a topic t . Further, coefficient $(1 - \alpha)$ is considered as biasing factor depending on the web page $P(v)$ selected for a queried topic t .

The page is having in-links from other pages. Assume N_t is number of topics to which a page is sensitive to surfing those topics. Effect of topics on PageRanks increases by using a non-uniform $N_t \times 1$ personalization vector for surfing probability p . Higher α means higher p .

Assume that the topics are t_1, t_2, \dots, t_n . Fix the number N_t . R_{TS} is to be computed for each of them. Therefore, compute for each topic t_j , the PageRank scores of page v as a function of t_j , which means compute $R(v, j)$, where $j = 1, 2, \dots, n$. That also means compute the n elements of a non-uniform $N_t \times 1$ personalization vector $R_{TS}(v)$ for t_1, t_2, \dots, t_n .

9.4.5 Hubs and Authorities

A hub is an index page that out-links to a number of content pages. A content page is topic authority. An authority is a page that has recognition due to its useful, reliable and significant information.

Figure 9.10(a) shows hubs (shaded circles) with the number of out-links associated with each hub. Figure 9.10(b) shows authorities (dotted circles) with the number of in-links and out-links associated with each link.

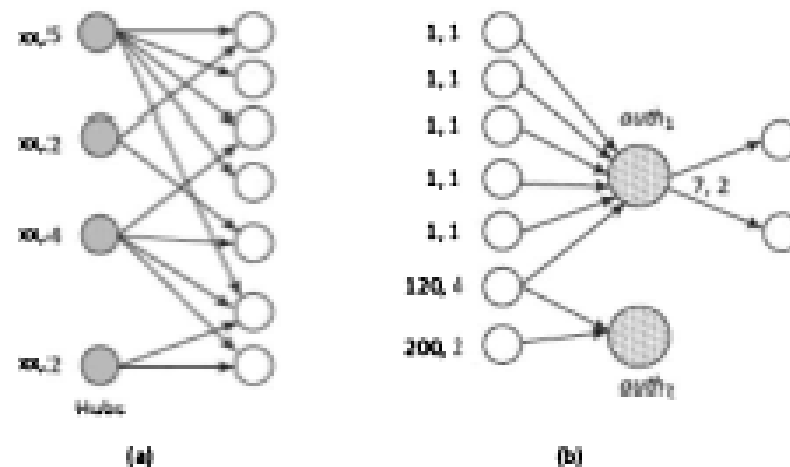


Figure 9.10 (a) Hubs (shaded circles) and (b) Authorities (dotted circles)

In-degrees (number of in-edges from other vertices) can be one of the measures for the authority. However, in-degrees do not distinguish between an in-link from a greater authority or lesser authority.

Authority, $auth_1$ in Figure 9.10(b) has in-links from 6 vertices (in-degrees = 6) and $auth_2$ has in-links to just 2 (in-degree = 2). However, $auth_1$ has link with six vertices with in-degrees = 1, 1, 1, 1, 1 and 120 (total = 125). Authority, $auth_2$ has links with two vertices with in-degrees = 120 and 200 (total = 220). $auth_2$ has association with greater authorities. Therefore, in-degrees may not be a good measure as compared to authority.

Kleinberg (1998) developed the Hypertext-Induced Topic Selection (HITS) algorithm.⁶ The algorithm computes the hubs and authorities on a specific topic t . The HITS analyses a sub-graph of web, which is relevant to t . Basis of computation is (i) hubs are the ones, which out-link to number of authorities, and (ii) authorities are the ones, which in-link to number of hubs. A bipartite graph exists for the hubs and authorities.

Consider a specifically queried topic t . Following are the steps:

1. Let a set of pages discover a root set R using standard search engine. Root pages may limit to top 200 for t .
2. Find a sub-graph of pages S , using a query that provides relevant pages for t and pointed by pages at R . Sub-graph S pages form Set for computations as it includes the children of parent R and limit to a random set of maximum 50 pages returned by a “reverse link” query.
3. Eliminate purely navigational links and links between two pages on the same host.
4. Consider only u ($|u| = 4-8$) pages from a given hyperlink as pointer to any individual page. (Section 9.4.2)

Sub-graph for HITS consisting of root set R of pages and children of parents in the sub-graph S . Figure 9.11 shows subgraph S for HITS consisting of root set R of pages and all the pages pointed to by any page of R .

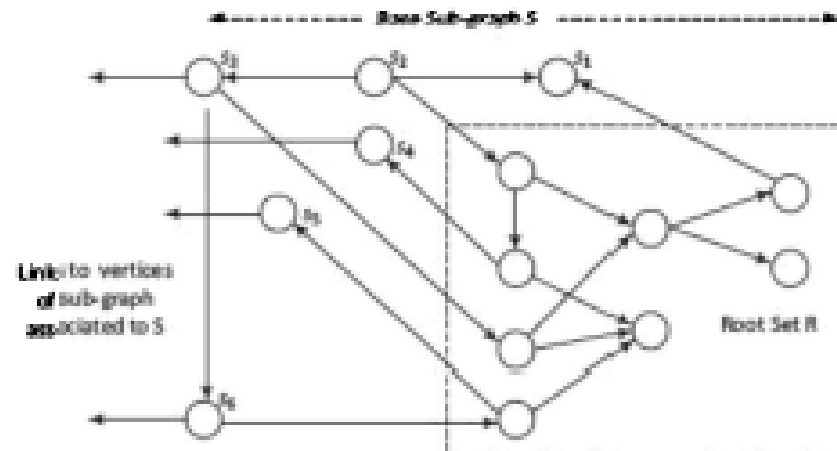


Figure 9.11 Sub-graph for HITS consisting of root set R of pages and base sub-graph S including all the pages pointed to by any page of R .

The left directed leftmost arrows from s_3 , s_4 , s_5 and s_6 are pointing to nodes in sub-graph(s) associated to S . The following example explains the algorithm steps to compute hub score and authority score.

9.4.6 Web Communities

Web communities are web sites or collections of websites, which limit the contents view and links to members. Examples of web communities are social networks, such as LinkedIn, SlideShare, Twitter and Facebook.

The communities consist of sites for do-it-yourself sites, social networks, blogs or bulletin boards. The issues are privacy and reliability of information.

Metric for analysis of web-community sites are web graph parameters, such as triangle count, clustering coefficient and K-neighbourhood.

K-neighbourhood analysis means the number of 1st neighbour nodes, 2nd neighbour nodes, and so on (K = 1, 2, 3, 4 and so on).

K-core analysis means the number of cores within a marked area. A core may consist of a triangle of connected vertices. A core may consist of a rectangle with interconnected edges and diagonals. A core may also be a group of cores.

Spark GraphX (Section 8.5) described functions for degree centralities, degree distribution, separation of degree, betweenness centralities, closeness centralities, neighbourhoods, strongly connected components, triangle counts, PageRank, shortest path, Breadth First Search (BFS), minimum spanning tree (forest), spectral clustering and cluster coefficient.

9.4.7 Limitations of Link, Rank and Web Graph Analysis

Following are the limitations of link and web graph analysis:

1. Search engines rely on metatags or metadata of the documents. That enhances the rank if metadata has biased information.
2. Search engines themselves may introduce bias while ranking the pages of clients higher as the pages of advertising companies may provide higher searches and hence lead to biased ranks.
3. A top authority may be a hub of pages on a different topic resulting in increased rank of the authority page.
4. Topic drift and content evolution can affect the rank. Off-topic pages may return the authorities.
5. Mutually reinforcing affiliates or affiliated pages/sites can enhance each other's rank and authorities.
6. The ranks may be unstable as adding additional nodes may have greater influence in rank changes.