

Analysis of Pentaho

BHAVESH REDDY MERUGUREDDY^{1,*}

¹School of Informatics and Computing, Bloomington, IN 47408, U.S.A.

* Corresponding authors: bmerugur@uemail.iu.edu

Paper 1, February 27, 2017

Pentaho is a leading business analytics and data integration tool that provides a qualified open source-based platform to assist a variety of big data deployments. It enables different organizations to utilize their data which helps them in delivering their services efficiently with minimum risk. Pentaho is often considered as an ideal application which can be used by businesses that desire to get the most out of their data and can also be used for embedded analytics.

© 2017 <https://creativecommons.org/licenses/>. The authors verify that the text is not plagiarized.

Keywords: Pentaho, Data Integration, Big Data, Community, ETL, MapReduce, SQL, Hadoop, OLAP

<https://github.com/bhaves37/sp17-i524/blob/master/paper1/S17-IR-2018/report.pdf>

INTRODUCTION

Pentaho can be viewed as a business intelligence suite that provides data mining, reporting, dashboarding and data integration capabilities. Generally, organizations tend to obtain meaningful relationships and useful information from the data present with them. Pentaho addresses the obstacles that obstruct them from doing so [1]. The platform includes a wide range of tools that analyze, explore, visualize and predict data easily which simplifies blending the data. The sole objective of pentaho is to translate data into value. Being an open and extensible source, pentaho provides big data tools to extract, prepare and blend any data. Along with this, the visualizations and analytics will help in changing the path that the organizations follow to run their business. From spark and hadoop to noSQL, pentaho transforms big data into big insights.

PENTAHO COMMUNITY

Pentaho provides two different editions, community edition and enterprise edition. As the name suggests, the enterprise edition comes with more packages to provide addition support. The community edition enables the developers or users to create complex solutions for the problems pertaining to their business [2]. The pentaho community has a group of intellectual people and helps the users in becoming a part of them and benefit from the open source contributions. These open source projects are helpful in delivering reliable and faster products which are timely tested by the community. The community includes all users like developers, testers and managers. Generally, the community edition platform enables the developers to sketch their design and develop a rough version of their product after which they can upgrade to enterprise edition for final production.

Pentaho provides an interactive console to its users. With

a few clicks of the mouse, users are allowed to interact with new data models and data. The platform hides the database connections and underlying application server and provides access to various data sources [3]. It provides metadata management capabilities and a dashboard to allow the administrators set security levels, monitor servers and set user access. There are many server plugins and desktop applications provided by pentaho.

Server applications

Business Intelligence platform is a basic service that provides reports, displays dashboards, reports business rules and performs OLAP analysis. The latest version comes with RESTful services and re-written scheduler along with a migration system. It generally runs in Apache java application server and can be embedded in any other java application server [1]. Pentaho analysis service is another server application that is written in java which primarily focuses on online analytical processing. It aggregates data into a memory cache by performing read operation from data sources like SQL. It comes with the pentaho platform in both the editions. These are some of the server applications provided by pentaho.

Desktop applications

Pentaho data mining is a desktop application that searches for patterns in data by performing knowledge analysis. All the techniques of data mining such as classification, clustering, regression and visualization are employed by this application along with some machine learning algorithms. This helps the users in predicting the trends in future. Pentaho metadata editor is an application that is used as an abstraction layer from the underlying data sources and helps the users in creating effective business models which can be used by other applications in

creating reports for the analytics. There are many more useful desktop applications.

Server plugins

Some of the important server plugins are community data access and data browser. Community data access is a pentaho server plugin that provides a common layer on the business analytics server for an easy data access. It runs the server by providing a REST interface and gets back the results in various forms such as xml, csv or json. Community data browser is a plugin that helps R in performing analytics on the data. It does the job of supplying queries to R by using online analytical processing browser.

DATA INTEGRATION

Extract, transform and load (ETL) are the basic operations that act as a tool for transforming data from one database and placing in other database. These processes can be carried out in pentaho with the help of a component called pentaho data integration, which is also referred as kettle [4]. The most useful functions of pentaho data integration include massive load of data into databases, data cleansing, migrating data between applications and integrating several applications. It is metadata oriented and can be used as a standalone application. The ability of transforming data is so high that the data can be manipulated with a very few limitations. Various input and output formats such as datasheets and text files are supported by pentaho data integration.

The transformation process undergoes three steps, input step, transformation step and output step. In the input step, data is ingested [5]. The data is then processed within pentaho data integration and the transformed data is given out in the output step. All these steps are carried out in parallel. The throughput of transformation process is restricted to speed of the step which is slowest. The slowest step is often referred as bottleneck. To improve the performance of transformation process, two steps are run in a loop which are, identification of the bottleneck and continued improvement of bottleneck until it is no longer a bottleneck.

Pentaho data integration has a set of components that contribute to its functionalities. They are spoon, kitchen, pan and carte [6]. Spoon can be considered as a desktop application that creates simple and even complex extract, transform and load (ETL) jobs without making the users write or read code. Spoon is the application that is used for transformations and jobs with the help of editor. So, it is the one that is used in most of the cases such as editing, debugging or running a transformation or a job. As the transformations are created in spoon, they can be executed with the help of a standalone command line process called pan. It is an engine that reads data, manipulates it and loads into various data sources. Kitchen is another standalone command line process that for executing jobs. It schedules different jobs to run at regular intervals. Carte provides remote execution capabilities and a medium for setting up a remote ETL server.

ARCHITECTURE

Pentaho architecture can be considered as a set of four components which are presentation layer, business intelligence platform, data and application integration and third party applications. Data can be provided to the presentation layer by reporting, analysis or process management. This data can then be

accessed through a web service, portal or a browser [7]. The security and repository issues are dealt by the business intelligence platform. Data integration and third party applications are respectively, the integration layer and applications with database from various sources.

The architecture also includes a set of predefined layers such as data layer, server layer and client layer. Data layer allows an application to connect to a data source. Server layer serves as a middle layer and several applications run on the server. Dashboards are provided to the end users by deploying them on the server along with the required reports. As mentioned above, a user console is provided that is used for security and configuration purposes. Client layer is of two forms, thin client and thick client. Thin client generally runs on a server. Analyzer and dashboard editor can be considered as the examples. Report designer and data integration come under thick client which act as a standalone.

BIG DATA USE CASES

Big data refers to humongous volumes of data being taken from multiple data sources and put into data stores. A use case can be defined as a technology solution for business specific challenges. Big data use cases help in understanding the problems that big data addresses.

Cyber security analysis helps the end users such as data scientists and security analysts in quickly detecting the threats. Cyber security analytics allows the users to utilize most of the staff resources via automation [8]. It empowers the data scientists with predictive analytics with the help of machine learning tools. It also provides the automation of blending and reporting on a variety of data. Pentaho platform can be utilized for data processing, data ingestion and delivery of threat calls with minimal costs and complexity.

Pentaho optimizes data warehouse and speeds up the development and deployment processes. It employs a simplified process for offloading to Hadoop. The offloaded data is usually less frequent data. Hand coding in MapReduce jobs and SQL can be avoided by the usage of visual integration tools. It provides access to data sources ranging from relational to operational to NoSQL technologies. Pentaho MapReduce helps in achieving high performance in a cluster environment. It provides a graphical and intuitive big data integration.

Another use case identified by pentaho is the streamlined data refinery. Pentaho data integration processes and refines different data sets by using Hadoop as its data processing platform. It provides modelled, delivered and published data sets to the users for visual analytics just by a mouse click. It can be seen as an integration process that blends huge volumes of highly diversified data. It also supplies tools for in-cluster simplified data processing and is regarded as a highly practical approach.

Pentaho's big data support extends the 360-degree view to internal and external customer related data. Customer service teams are provided with time-sensitive and blended streams of data. This helps in making profitable decisions. The presence of an adaptive big data layer relieves several organizations from evolving technologies. Customers are given access to customizable, intuitive and interactive dashboards. Data scientists are provided with predictive analytics and data mining tools.

Monetizing the data is the final use case addressed by pentaho. It allows the users to capitalize on big data with the help of powerful data processing and embeddable data analytics [9]. Pentaho's big data analytics platform empowers easy big data

ingestion and transformation as it works as a no-code data integration environment. It is a flexible platform that supports security and deployments specific to customers.

COMPARISON

Pentaho products compete with some big names in current field such as SAP, IBM and oracle. Pentaho provides open source solutions and is considered to be much cheaper than the proprietary equivalents. Jaspersoft is an established open source rival of pentaho. Though both pentaho and jaspersoft offer similar features with similar costs, pentaho has got wider online presence and more followers in social media [10].

CONCLUSION

Pentaho is an open source based platform for diverse big data deployments. It empowers analytics in any environment by delivering governed data. It has unified data integration and analytics components which are comprehensive and embeddable. The primary aim of pentaho is to enable organizations to find new revenue streams with extraordinary service at minimum risk. It helps them in harnessing the value from their data in order to make their operations efficient and consistent.

REFERENCES

- [1] "Pentaho," webpage. [Online]. Available: <https://en.wikipedia.org/wiki/Pentaho>
- [2] "Community wiki home," Webpage. [Online]. Available: <http://wiki.pentaho.com/display/COM/Community+Wiki+Home>
- [3] "Pentaho bi suite enterprise edition," Webpage, 2006. [Online]. Available: <http://searchdatamanagement.techtarget.com/review/Pentaho-BI-Suite-Enterprise-Edition>
- [4] M. C. Roldán, "Pentaho data integration (kettle) tutorial," Webpage, 2008. [Online]. Available: [http://wiki.pentaho.com/display/EAI/Pentaho+Data+Integration+\(Kettle\)+Tutorial](http://wiki.pentaho.com/display/EAI/Pentaho+Data+Integration+(Kettle)+Tutorial)
- [5] R. Haces, "Pentaho data integration performance tuning," Webpage. [Online]. Available: <https://support.pentaho.com/hc/en-us/articles/205715046-Best-Practice-Pentaho-Data-Integration-%20Performance-Tuning->
- [6] "Pentaho data integration architecture," Webpage. [Online]. Available: <https://help.pentaho.com/Documentation/5.3/OL0/0Y0/010>
- [7] "Understanding pentaho architecture," Webpage. [Online]. Available: <https://www.edureka.co/blog/understanding-pentaho-architecture/>
- [8] "What is big data?" Webpage. [Online]. Available: <http://www.pentaho.com/what-is-big-data#tab-3>
- [9] "Monetize my data," Webpage. [Online]. Available: <http://www.pentaho.com/Monetize-My-Data>
- [10] "Compare pentaho vs. jaspersoft," Webpage. [Online]. Available: <https://comparisons.financesonline.com/jaspersoft-vs-pentaho>