



CENTRO DE INFORMÁTICA  
GRADUAÇÃO

BRUNO DE ASSIS PEREIRA

**Um Processo para Desenvolvimento de um *Data Warehouse* em um cenário de *Big Data* utilizando a ferramenta *Hadoop***

TRABALHO DE GRADUAÇÃO

Recife  
2016

BRUNO DE ASSIS PEREIRA

**Um Processo para Desenvolvimento de um *Data Warehouse* em um cenário de *Big Data* utilizando a ferramenta *Hadoop***

Trabalho de conclusão de curso apresentado à disciplina de TG como parte dos requisitos necessários à obtenção do título de Bacharel em Ciências da Computação.

Orientador: Fernando da Fonseca de Souza

Recife  
2016

## **Agradecimentos**

Agradeço a todos que me acompanharam durante todo o desenvolvimento deste trabalho e também agradeço a todos que estiveram ao meu lado durante todo o curso de graduação seja me ensinando, trabalhando ou aprendendo comigo. Agradeço principalmente aos meus pais, Marisa Assis e Leonel Pereira, aos meus irmãos, Rodrigo Pereira e Gabriela Khuni, a minha namorada Thais Frota, e aos meus primos e amigos. Agradeço também ao meu orientador Fernando Fonseca por ser um excelente professor ao me guiar não só neste trabalho, mas no decorrer do curso.

## Resumo

Dados são a codificação das informações e do conhecimento. Realizar um gerenciamento sobre essa informação envolve um processo de organizar, adquirir, armazenar, recuperar, e gerenciar esses dados. Estes são coletados por meio de diferentes processos e usados para auxiliar tomadas de decisão, de modo que aqueles que estão executando e consumindo os resultados do processo possam compreender toda a informação e atender aos seus diversos requisitos.

Um elemento que evoluiu para dar suporte ao processo de tomada de decisões foi o *Data Warehouse* (DW) (SUN et al., 2014). Por meio da capacidade de coletar, armazenar, e gerenciar dados, aplicando métodos tradicionais e estatísticos de medições para criar relatórios e plataformas de análise, o *Data Warehousing* se tornou um elemento chave no processo de tomada de decisões. Esta última pode se tornar ainda mais segura caso se possa analisar um número maior de dados.

O *Big Data* (MELOROSE; PERROY; CAREAS, 2015), conjunto de dados que se caracteriza pelo seu grande volume gerado em alta velocidade e em diversos tipos e formatos, tem se popularizado principalmente pela seguinte razão: as plataformas tecnológicas que surgiram junto com esse complexo e amplo conjunto de dados, provêm a capacidade de processar vários formatos e estruturas de dados sem ter que se preocupar com as restrições associadas aos sistemas de bancos de dados tradicionais. Como exemplo de ferramentas tem-se o *Hadoop* (JURNEY, 2013) que é uma arquitetura proposta como solução para o processamento de *Big Data* em uma plataforma mais barata com rápida escalabilidade e processamento paralelo.

Construir um *Data Warehouse* que contém *Big Data* possibilitará o desenvolvimento de análises sobre esses dados, e esta atividade trará conhecimento valioso que antes estava escondido ou não estaria disponível tão facilmente.

Neste trabalho serão explorados projetos que realizaram análises sobre *Big Data*. Ao fim dos estudos, a construção de um DW sobre um cenário de *Big Data* utilizando o *Hadoop* será realizada a fim de apresentar por meio da criação de relatórios contendo aplicação de métodos estatísticos, os conceitos que foram extraídos e o potencial que esse tipo de desenvolvimento tem.

## Lista de ilustrações

Figura 1 – Estrutura do DW . . . . .	14
Figura 2 – Os três V's do Big Data . . . . .	15
Figura 3 – Componentes do DW . . . . .	16
Figura 4 – Processamento de dados tradicional . . . . .	18
Figura 5 – Processamento de dados <i>Big Data</i> . . . . .	19
Figura 6 – Data Science no contexto de vários processos relacionados aos dados de uma organização . . . . .	20
Figura 7 – Estrutura do HDFS . . . . .	23
Figura 8 – Estrutura do Hive . . . . .	26
Figura 9 – Arquitetura do Data Warehousing . . . . .	40
Figura 10 – Interfaces da Cloudera . . . . .	41
Figura 11 – Diagrama de Classes . . . . .	42
Figura 12 – Arquitetura do DW - TPCx-BB & Hive . . . . .	59

## **Lista de quadros**

Quadro 1 – Critério de Análise e Apresentação . . . . .	30
Quadro 2 – Sumário de Análises . . . . .	38

## Lista de tabelas

Tabela 1 – Dimensão de Datas. . . . .	43
Tabela 2 – Dimensão de Clientes. . . . .	44
Tabela 3 – Endereço do cliente . . . . .	45
Tabela 4 – Dados demográficos do cliente . . . . .	46
Tabela 5 – Dados demográficos do domicílio . . . . .	47
Tabela 6 – Inventário dos itens estocados no warehouse . . . . .	47
Tabela 7 – Dimensão Item . . . . .	48
Tabela 8 – Dimensão com informações sobre revisões dos produtos . . . . .	48
Tabela 9 – Dimensão com informações sobre as promoções . . . . .	49
Tabela 10 – Dimensão com informações sobre as razões de retorno dos itens .	50
Tabela 11 – Dimensão de lojas . . . . .	50
Tabela 12 – Dimensão de tempo . . . . .	52
Tabela 13 – Dimensão de armazéns . . . . .	52
Tabela 14 – Dimensão de Página da Web . . . . .	53
Tabela 15 – Dimensão de Web Site . . . . .	54
Tabela 16 – Fato Preços de Mercado dos Itens . . . . .	55
Tabela 17 – Fato Retornos de Vendas . . . . .	55
Tabela 18 – Fato Vendas da Loja . . . . .	57

# Sumário

<b>1</b>	<b>Introdução</b>	<b>10</b>
1.1	Motivação	10
1.2	Objetivos	11
<b>1.2.1</b>	<b>Objetivos Específicos</b>	<b>11</b>
1.3	Metodologia	11
1.4	Estrutura do Trabalho	12
<b>2</b>	<b>Fundamentação Conceitual</b>	<b>13</b>
2.1	Data Warehousing	13
2.2	Data Warehousing e Big Data	14
2.3	Big Data	17
<b>2.3.1</b>	<b>Tecnologias de <i>Big Data</i></b>	<b>19</b>
<b>2.3.2</b>	<b>Hadoop</b>	<b>20</b>
<b>2.3.3</b>	<b>HDFS</b>	<b>21</b>
<b>2.3.4</b>	<b>MapReduce</b>	<b>23</b>
<b>2.3.5</b>	<b>Hive</b>	<b>23</b>
2.3.5.1	Arquitetura do Hive	24
<b>2.3.6</b>	<b>Avro</b>	<b>26</b>
<b>2.3.7</b>	<b>Cassandra</b>	<b>27</b>
<b>2.3.8</b>	<b>Chuckwa</b>	<b>27</b>
<b>2.3.9</b>	<b>HBase</b>	<b>27</b>
<b>2.3.10</b>	<b>Mahout</b>	<b>28</b>
<b>2.3.11</b>	<b>Pig</b>	<b>28</b>
<b>2.3.12</b>	<b>ZooKeeper</b>	<b>28</b>
2.4	Comentários Finais	28
<b>3</b>	<b>Trabalhos Relacionados</b>	<b>30</b>
3.1	Trabalho 1 – Análises sobre Imagens Médicas (ZHANG et al., 2016)	30
<b>3.1.1</b>	<b>Problema</b>	<b>30</b>
<b>3.1.2</b>	<b>Solução</b>	<b>31</b>
<b>3.1.3</b>	<b>Limitações</b>	<b>31</b>
3.2	Trabalho 2 – Transformando o cenário do Mercado com o Cloudera	31
<b>3.2.1</b>	<b>Problema</b>	<b>31</b>
<b>3.2.2</b>	<b>Solução</b>	<b>32</b>
<b>3.2.3</b>	<b>Limitações</b>	<b>32</b>

3.3	Trabalho 3 – Aprimorando a conectividade de assistências médica com <i>Big Data</i> . . . . .	33
3.3.1	<b>Problema</b> . . . . .	33
3.3.2	<b>Solução</b> . . . . .	33
3.3.3	<b>Limitações</b> . . . . .	33
3.4	Trabalho 4 – Um Sistema Analítico de <i>Big Data</i> para analisar a percepção de segurança dos cidadãos (CAMARGO et al., 2016) . . . . .	34
3.4.1	<b>Problema</b> . . . . .	34
3.4.2	<b>Solução</b> . . . . .	34
3.4.3	<b>Limitações</b> . . . . .	35
3.5	Trabalho 5 – Aplicando transformações centradas no cliente . . . . .	36
3.5.1	<b>Problema</b> . . . . .	36
3.5.2	<b>Solução</b> . . . . .	36
3.5.3	<b>Limitações</b> . . . . .	36
3.6	Comentários Finais . . . . .	37
<b>4</b>	<b>Arquitetura e Processo de Desenvolvimento do DW</b> . . . . .	<b>40</b>
4.1	Arquitetura do DW . . . . .	40
4.2	Especificação das Ferramentas . . . . .	41
<b>4.2.1</b>	<b>Desenvolvimento do DW</b> . . . . .	<b>42</b>
4.2.1.1	Dimensões . . . . .	43
4.2.1.1.1	<i>Data_dim</i> . . . . .	43
4.2.1.1.2	<i>Cliente</i> . . . . .	44
4.2.1.1.3	<i>Endereco_cliente</i> . . . . .	45
4.2.1.1.4	<i>Demografia_cliente</i> . . . . .	46
4.2.1.1.5	<i>Demografia_domicilio</i> . . . . .	47
4.2.1.1.6	<i>Inventario</i> . . . . .	47
4.2.1.1.7	<i>Item</i> . . . . .	47
4.2.1.1.8	<i>Revisoes_Produtos</i> . . . . .	48
4.2.1.1.9	<i>Promocao</i> . . . . .	49
4.2.1.1.10	<i>Razao</i> . . . . .	50
4.2.1.1.11	<i>Loja</i> . . . . .	50
4.2.1.1.12	<i>Tempo_dim</i> . . . . .	51
4.2.1.1.13	<i>Armazem</i> . . . . .	52
4.2.1.1.14	<i>Pagina_Web</i> . . . . .	53
4.2.1.1.15	<i>Web_site</i> . . . . .	53
4.2.1.2	<b>Fatos</b> . . . . .	55
4.2.1.2.1	<i>Item_precos_mercado</i> . . . . .	55
4.2.1.2.2	<i>Retorno_Vendas</i> . . . . .	55
4.2.1.2.3	<i>Vendas_Loja</i> . . . . .	57

<b>4.2.2</b>	<b>Volume de Dados</b> . . . . .	<b>58</b>
<b>4.2.3</b>	<b>Processo de ETL</b> . . . . .	<b>58</b>
4.3	Comentários Finais . . . . .	59
<b>5</b>	<b>Conclusão</b> . . . . .	<b>60</b>
5.1	Limitações . . . . .	60
5.2	Trabalhos Futuros . . . . .	60
	<b>Referências</b> . . . . .	<b>62</b>

# 1 Introdução

A Internet mudou e continua a causar grande impacto sobre o comportamento da sociedade, pois essa grande disseminadora de informações altera a forma como são gerenciados os negócios, a interação com os clientes, o compartilhamento de informações entre pessoas, a determinação do valor de um determinado produto, e, além disso, a Web proveu um novo canal para geração de lucro. Quando se deseja receber informação a respeito de novos produtos, ou até mesmo entender sobre as dificuldades já vivenciadas por outras pessoas ao utilizar serviços e produtos de uma determinada empresa, realiza-se uma busca na Internet.

No entanto, não são somente os consumidores que possuem interesse na informação disponibilizada no mundo das redes. Buscando entender melhor as necessidades de seus clientes, as pequenas e grandes empresas utilizam dos dados gerados pelos seus *stakeholders*<sup>1</sup> para realizar de forma mais precisa o processo de tomada de decisão.

A principal função da aplicação de análises sobre os dados é de realizar tomadas de decisão, e para auxiliar nessa tarefa criou-se o Data Warehouse (DW). O Data Warehousing foi desenvolvido para dar suporte ao processo de tomada de decisões por meio da sua capacidade de coletar, armazenar, e gerenciar dados, aplicando métodos tradicionais e estatísticos de medições para criar relatórios e plataformas de análise.

Atualmente, vive-se a era do *Big Data* na qual grandes volumes de dados de diversos formatos são gerados em alta velocidade. Construir um *Data Warehouse* contém *Big Data* possibilitará o desenvolvimento de análises sobre esses dados, e esta atividade trará conhecimento valioso que antes estava escondido ou não estaria disponível tão facilmente.

Uma das ferramentas mais utilizadas para resolver o processamento de *Big Data* é o *Hadoop* que possui uma rápida escalabilidade e utiliza de processamento paralelo. Além disso, essa ferramenta também faz uso de plataformas mais baratas tornando-a assim mais acessível.

## 1.1 Motivação

Essa grande disponibilidade dos dados tem levado a um crescimento no interesse sobre os métodos para extrair conhecimento e informação útil dos mesmos – *Data*

---

<sup>1</sup> pessoas, grupos ou organizações que tem interesse ou se importam com uma organização; pessoas, grupos ou organizações que podem afetar ou podem ser afetados por ações, objetivos e políticas de uma organização.

*Science* (Ciência dos dados) (PROVOST, 2013). O entendimento sobre as ferramentas que fazem parte dessa área possibilita que se possa fazer parte desse novo mundo rico em oportunidades para realizar grandes descobertas. Conseguir desenvolver uma plataforma com a qual o processo de tomada de decisão se torna ainda mais eficiente traz benefícios para diversas áreas – educação, saúde, aviação, desenvolvimento sustentável, construção, dentre outras.

## 1.2 Objetivos

Este trabalho tem como objetivo principal analisar e compreender mais sobre o processo de análises sobre Big Data, utilizando técnicas de DW.

### 1.2.1 Objetivos Específicos

Para alcançar o objetivo principal, são definidos os seguintes objetivos específicos:

- Desenvolver um processo de criação de um Data Warehouse que envolva Big Data;
- Construir o DW utilizando o Hadoop;
- Criar relatórios contendo resultados da aplicação de métodos estatísticos, os conceitos que foram extraídos; e
- Destacar o potencial existente na aplicação de técnicas de Data Warehousing sobre Big Data.

## 1.3 Metodologia

Com o objetivo de proporcionar um maior aprendizado sobre o tema de *Data Warehousing* aplicado a *Big Data* uma metodologia exploratória foi aplicada. Com este tipo de pesquisa, é possível se familiarizar com o fenômeno investigado de modo que o próximo passo da pesquisa possa ser melhor compreendido e com maior precisão (EDUCAÇÃO, 2013).

Para realizar este trabalho, inicialmente foi realizado um estudo sobre o estado da arte de modo a entender sobre o atual cenário *Big Data* e como aplicar conceitos de *Data Warehouse* ao mesmo. Com isso, foi encontrada a ferramenta *Hadoop*. Após realizar estudos sobre a ferramenta, com a utilização da máquina virtual da *Cloudera Distribution for Hadoop* (CDH)<sup>2</sup> e do *framework* de geração de dados da TPC (*Transac-*

<sup>2</sup> <http://www.cloudera.com/downloads/cdh/5-8-2.html>

tion Processing Performance Council) *Express Big Bench* (TPCx-BB)<sup>3</sup> foi construído um DW para os dados gerados com o *Hive*<sup>4</sup>. Com esse DW foi possível extrair informações e apresentá-las de modo a mostrar o potencial de um DW sobre um possível cenário *Big Data*.

A fins de estudo e compreensão das tecnologias, foi utilizado o Hadoop no modo pseudo distribuído (*pseudo distributed mode*), pois este é recomendado para testes pelo fato de não ser necessário utilizar um *cluster* de máquinas e conseguir simular um cluster com apenas uma máquina. Dessa forma, para os objetivos levantados o modo escolhido já será suficiente para compreensão das ferramentas e apresentação de potencial.

#### 1.4 Estrutura do Trabalho

Este trabalho será dividido em cinco (5) capítulos. Além do capítulo de Introdução, tem-se a seguinte estrutura: No Capítulo 2 é apresentado o estado atual das técnicas de *Data Warehousing* sobre *Big Data*, de modo a trazer conceitos e características importantes sobre os dois temas. Além disso, também é apresentado o *Hadoop*, uma das maiores tecnologias de *Big Data*, fazendo uma descrição das principais funcionalidades da ferramenta e dissertando sobre as tecnologias que compõem o ecossistema do *Hadoop*, dando ênfase ao *Hive*.

No Capítulo 3 são analisados trabalhos relacionados, trazendo uma descrição dos mesmos, com os resultados obtidos e técnicas aplicadas.

No Capítulo 4 é apresentado o processo e a especificação das ferramentas utilizadas para o desenvolvimento do DW, e é descrita a implementação realizada.

No Capítulo 5 são apresentadas as contribuições e limitações do trabalho, bem como sugestões de possíveis trabalhos futuros.

Por fim, são listadas as referências bibliográficas utilizadas.

<sup>3</sup> <http://www.tpc.org/tpcx-bb/default.asp>

<sup>4</sup> <http://hive.apache.org/>

## 2 Fundamentação Conceitual

Neste capítulo serão descritos conceitos importantes para a compreensão do trabalho. Dentre estes conceitos, serão abordados temas como *Big Data*, as tecnologias do *Big Data* – o *Hadoop* e o seu ecossistema –, e *Data Warehousing*.

### 2.1 Data Warehousing

Os primeiros conceitos sobre *Data Warehouse* surgiram com a necessidade de armazenar e analisar os dados do OLTP<sup>1</sup>. A habilidade de reunir transações, produtos, serviços, e localizações durante um período de tempo começou a prover capacidades interessantes para as companhias que nunca foi possível no mundo do OLTP, devido ao *design* desses sistemas e devido às limitações da sua infraestrutura em termos de escalabilidade (MELOROSE; PERROY; CAREAS, 2015).

A definição de *Data Warehouse* que é aceita como padrão pela indústria afirma que o DW é uma coleção de dados orientada por assunto<sup>2</sup>, não volátil<sup>3</sup>, integrada<sup>4</sup>, e que varia no tempo<sup>5</sup>, criada para dar suporte sobre o processo de tomadas de decisão (INMON; STRAUSS; NEUSHLOSS, ).

Os elementos básicos que compõem a arquitetura de um *Data Warehouse* (Figura 1) são as seguintes:

- Fonte de dados<sup>6</sup>;
- ETL<sup>7</sup>;
- *Staging Area*<sup>8</sup>;

<sup>1</sup> Processamento de transações em tempo real; sistemas que se encarregam de registrar todas as transações contidas em uma determinada operação organizacional.

<sup>2</sup> Os sistemas transacionais são organizados de acordo com os principais assuntos da empresa em questão.

<sup>3</sup> No ambiente do DW, os dados, antes de serem carregados, são filtrados e limpos. Após esta etapa essas dados sofrem somente operações de consulta e exclusão, sem que possam ser alterados.

<sup>4</sup> A integração dos dados é realizada visando padronizar os dados dos diversos sistemas em uma única representação.

<sup>5</sup> Consiste na manutenção de um histórico de dados em relação ao período de tempo.

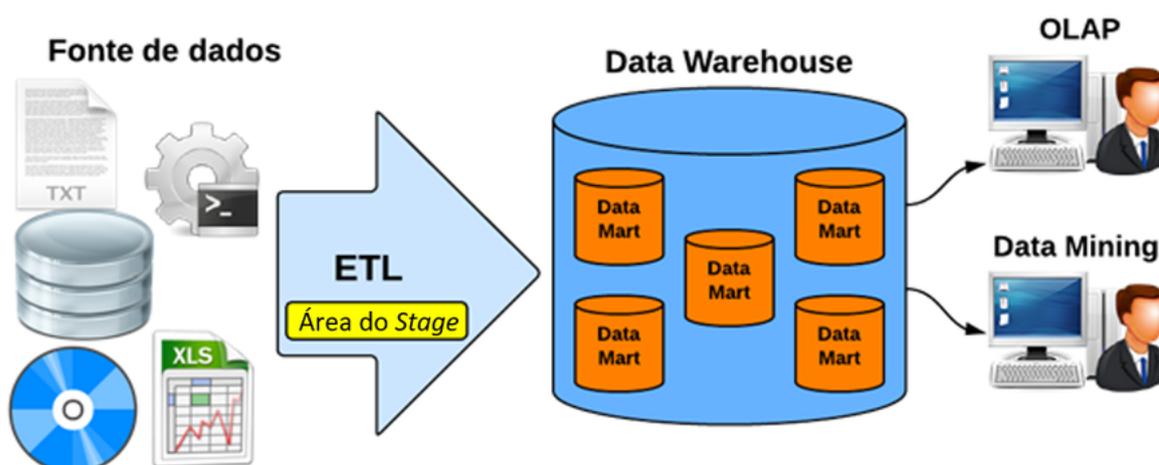
<sup>6</sup> Sistemas transacionais da empresa, pode ser composto por diversas formas de dados.

<sup>7</sup> Do Inglês *Extract, Transform and Load*, é o principal processo de condução dos dados até o armazenamento definitivo no DW. É responsável por todas as tarefas de extração, tratamento e limpeza de dados, e inserção na base do DW.

<sup>8</sup> Área de armazenamento intermediária situada dentro do processo de ETL. Auxilia a transição dos dados das origens para o destino final no DW.

- *Data Warehouse*<sup>9</sup>;
- *Data Mart*<sup>10</sup>;
- *OLAP*<sup>11</sup>; e
- *Data Mining*<sup>12</sup>.

Figura 1 – Estrutura do DW



Fonte: <https://corporate.canaltech.com.br/materia/business-intelligence/conhecendo-a-arquitetura-de-d-ata-warehouse-19266/>

Com o advento do *Data Warehouse*, as empresas podem atenciar a concorrência no que tange a gestão de informações que fornecem competitividade e inteligência no mercado, além de favorecerem o crescimento e alcance de resultados positivos na organização.

## 2.2 Data Warehousing e Big Data

As novas tecnologias que surgem estão muitas vezes ligadas à inovação, e essa última mudou a forma que se entra no mundo dos negócios e são providos serviços. Além disso, essas criações inovadoras alteram a forma de avaliação de valor associada a produtos e serviços, e as formas de lucro.

<sup>9</sup> Estrutura de armazenamento das informações decisivas. Apenas os dados com valor para a gestão corporativa estarão no DW.

<sup>10</sup> Estrutura similar ao DW, porém com uma proporção menor de informações.

<sup>11</sup> Do Inglês On-Line Analytical Processing, na arquitetura de um DW se refere às ferramentas com capacidade de análise em múltiplas perspectivas das informações armazenadas.

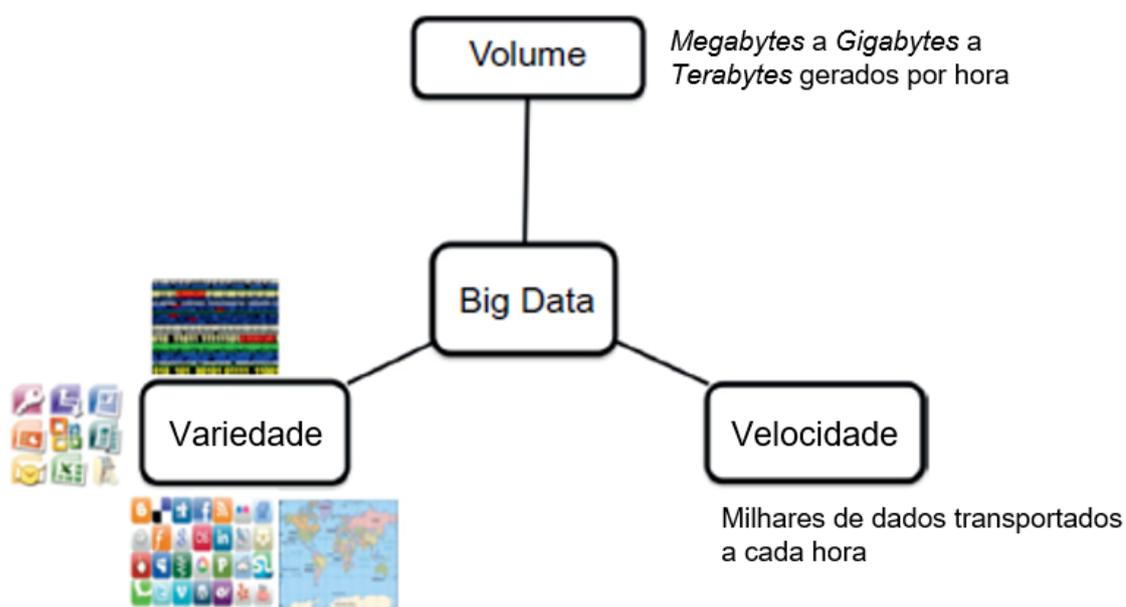
<sup>12</sup> A mineração de dados se refere às ferramentas com capacidade de descoberta de conhecimento relevante dentro do DW.

Essas tendências adicionaram complexidade em termos de processos, e ao mesmo tempo criaram a necessidade de adquirir dados necessários para esses processos, os quais podem prover conhecimentos essenciais para áreas nas quais nunca antes foi possível (MELOROSE; PERROY; CAREAS, 2015).

Com a vasta quantidade de dados agora disponível, as empresas em praticamente todas as áreas estão focadas em explorar essas informações visando obter vantagens competitivas no mercado (PROVOST, 2013).

Os desafios do crescimento de dados e oportunidades são definidos como tridimensionais (MELOROSE; PERROY; CAREAS, 2015). São caracterizados pelo volume<sup>13</sup> crescente, velocidade<sup>14</sup>, e variedade<sup>15</sup> - os três V's (Figura 2). A indústria utiliza essa definição como um padrão para classificar *Big Data* (BETTER et al., 2013).

Figura 2 – Os três V's do Big Data



Fonte: Livro "Data Warehousing in the Age of Big Data"

Para conseguir trabalhar com esse tipo de dados, a arquitetura dos *Data Warehouses* precisou sofrer adaptações. A nova geração de DW é mais complexa em se tratando do desenvolvimento de uma arquitetura física, consistindo de várias tecnologias, e vai ser orientada a dados numa perspectiva de integração de todos os dados de

<sup>13</sup> quantidade de dados disponíveis e gerados continuamente.

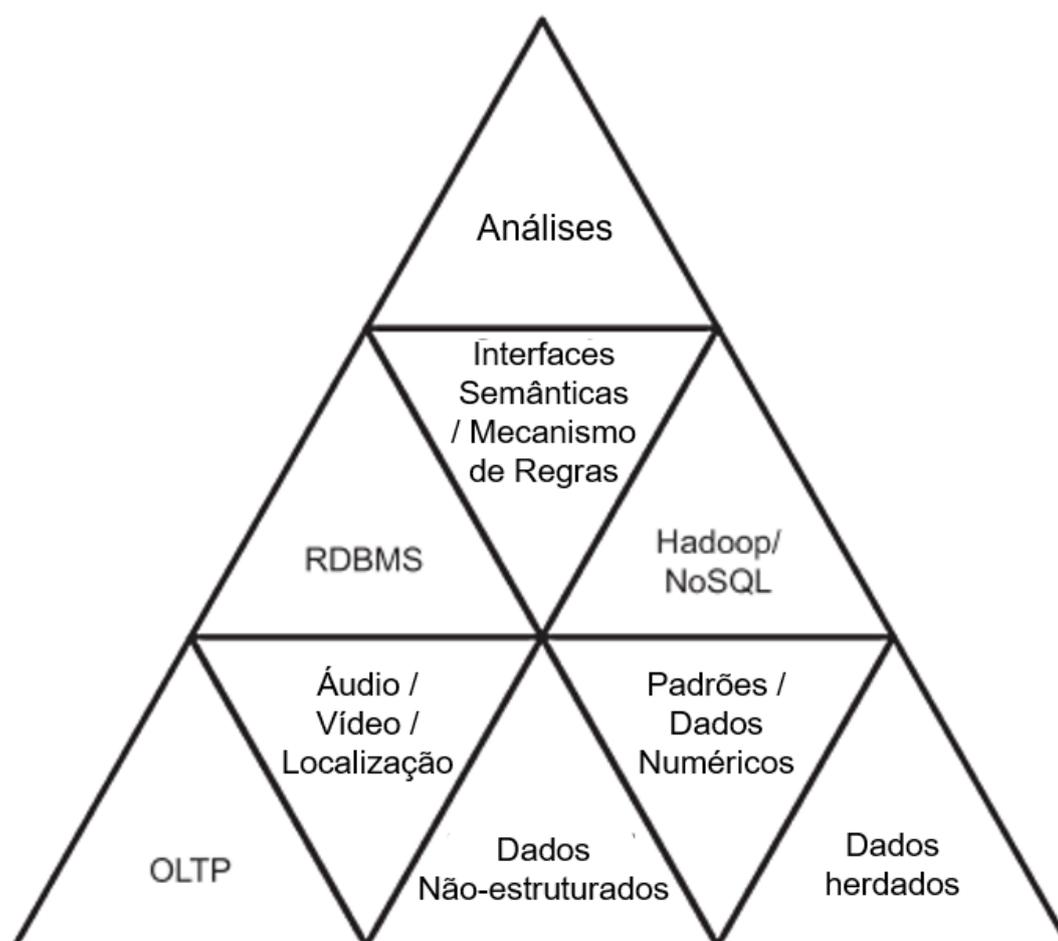
<sup>14</sup> velocidade com que os dados entram e saem, os dados são transmitidos de forma contínua e os conjuntos resultantes são úteis quando a aquisição e os atrasos de processamento levam pouco tempo.

<sup>15</sup> a extensão dos tipos de dados gerados, estes vem em diversos formatos que variam de e-mails a tweets, de rede social a dados de sensores. Não se tem um controle sobre o formato dos dados de entrada ou sobre a estrutura dos dados.

uma empresa que também terão *Big Data* como parte desse conjunto de informações. Além disso, essa geração será extremamente flexível e escalável em se tratando de uma perspectiva de arquitetura de dados (MELOROSE; PERROY; CAREAS, 2015).

Na Figura 3 podemos observar alguns dos componentes da próxima geração de *Data Warehouses*.

Figura 3 – Componentes do DW



Fonte: Livro “Data Warehousing in the Age of Big Data”

A primeira camada representa os dados, a segunda camada representa as tecnologias que vão ser utilizadas para a integração dos dados de vários tipos e fontes de dados, e a camada do topo representa a camada de análise que vai ser utilizada para orientar as necessidades da visualização da próxima geração de *business intelligence*<sup>16</sup> e análises (MELOROSE; PERROY; CAREAS, 2015).

<sup>16</sup> Pode ser descrito como um conjunto de técnicas e ferramentas utilizados para a aquisição e transformação de dados em informação útil e significativa para realização de análises de negócios.

A próxima geração das atividades referentes a arquitetura e processamento de dados precisará incluir algoritmos como, *Text Mining*<sup>17</sup>, *Data Mining*, *Pattern Processing*<sup>18</sup>, *Statistical Models*<sup>19</sup> e *Mathematical Models*<sup>20</sup>, para resolver problemas específicos de processamento de dados, criação de relatórios e análises a partir do DW sobre cenários com *Big Data* (MELOROSE; PERROY; CAREAS, 2015).

### 2.3 Big Data

O termo Big Data identifica conjuntos de dados de tipos específicos, principalmente de dados não-estruturados, que povoam a camada de dados das aplicações científicas de computação e da Web. Esses dados têm características específicas em comum:

- Dados em larga escala, que se refere ao tamanho e a distribuição dos repositórios de dados;
- Problemas de escalabilidade, que se refere às capacidades das aplicações ao executar repositórios de dados enormes e de larga escala (i.e., *big data*, em resumo) para escalar rapidamente sobre entradas de tamanho crescente;
- Suporte avançado aos processos de ETL<sup>21</sup> a partir de dados de baixo nível, dados brutos, para torná-los uma informação de certa forma estruturada; e
- Criação e desenvolvimento sobre análises fáceis e interpretáveis sobre repositórios de *Big Data* de modo a obter inteligência e extrair conhecimento útil destes (CUZZOCREA; SACCÀ; ULLMAN, 2013).

*Big Data* também é o território no qual os tradicionais bancos de dados relacionais e sistemas de arquivos têm suas capacidades de processamento excedidas por altos volumes de transação, pela velocidade de resposta, e pela quantidade e ou variedade de dados. Essa gigantesca quantidade de dados cobre uma variedade de situações, todas elas com a palavra “mais” em comum – mais variedade, mais quantidade, mais usuários, mais velocidade, mais complexidade (BETTER et al., 2013).

<sup>17</sup> Processamento de texto baseado em regras de negócio definidas pelo usuário a fim de extrair dados que podem ser usados na classificação do texto para futuras explorações sobre os dados.

<sup>18</sup> Ramo do aprendizado de máquina que foca no reconhecimento de padrões e regularidades nos dados.

<sup>19</sup> Uma classe do modelo matemático, que inclui um conjunto de suposições sobre a geração de amostras de dados.

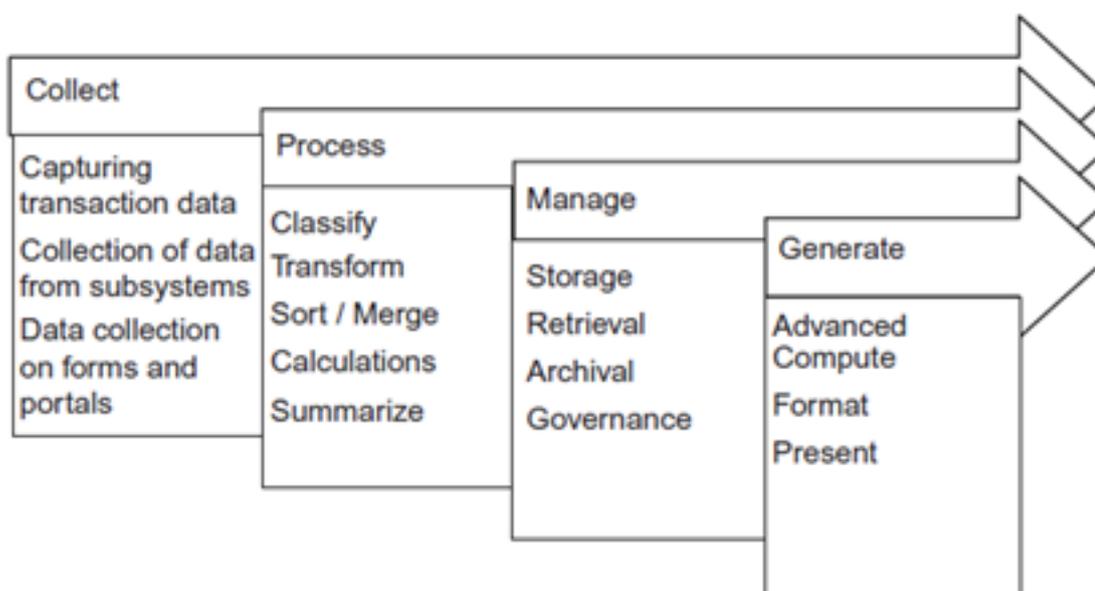
<sup>20</sup> Descrição de um sistema que utiliza conceitos e linguagens matemáticas.

<sup>21</sup> processo de extração de dados dos sistemas fonte para trazê-los para o data warehouse; Extrair, Transformar, e carregar (Load) dados.

A forma como se lida com o processamento de *Big Data* se dá por meio de processamento e armazenamento distribuídos, redes neurais, arquiteturas multiprocessadoras, e conceitos de orientação a objetos, combinados com técnicas de processamento de dados da Internet.

O processamento de dados pode ser definido como coleta, processamento, e gerenciamento de dados resultando em geração de informação para os usuários finais (Figura 4).

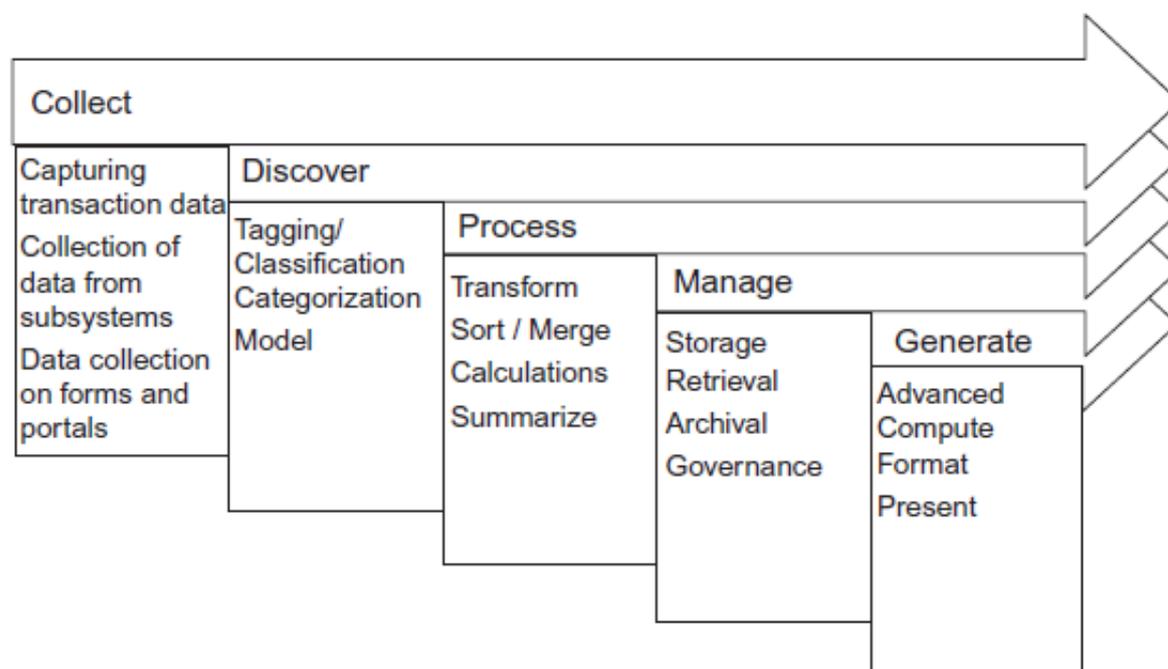
**Figura 4 – Processamento de dados tradicional**



Fonte: Livro “Data Warehousing in the Age of Big Data”

No entanto, o processamento de dados em se tratando de *Big Data* ocorre de forma diferente (Figura 5). O primeiro estágio é o de coleta de dados. Nesse estágio, os dados são recebidos a partir de diferentes fontes e carregados para um sistema de arquivos chamado de *landing zone*<sup>22</sup> ou *landing area*. Os dados são distribuídos em subdiretórios de acordo com o seu tipo. No estágio seguinte, os dados são carregados com a aplicação de metadados (esse é o momento em que pela primeira vez se aplica uma estrutura para os dados) e estes são preparados para a transformação. Na etapa de transformação, os dados são modificados por meio da aplicação das regras de negócio e do processamento do conteúdo. Por fim, o conjunto de dados resultantes pode ser extraído para processamentos futuros que incluem análises, relatórios operacionais, integração com *data warehouse*, e visualizações (MELOROSE; PERROY; CAREAS, 2015).

<sup>22</sup> centro onde os dados vão ficar armazenados; sistema de arquivos.

Figura 5 – Processamento de dados *Big Data*

Fonte: Livro “Data Warehousing in the Age of Big Data”

### 2.3.1 Tecnologias de *Big Data*

Os sistemas de processamento de dados tradicionais não têm capacidade para realizar tal tarefa sobre esses conjuntos de dados que são muito grandes, por esse motivo são necessárias novas tecnologias. Para realizar o processamento de *Big Data*, tecnologias como *Hadoop*<sup>23</sup>, *Hive*<sup>24</sup>, *HBase*<sup>25</sup>, *MongoDB*<sup>26</sup> começaram a ser amplamente utilizadas (MELOROSE; PERROY; CAREAS, 2015).

Assim como as tecnologias tradicionais, as tecnologias de *Big Data* são utilizadas para várias tarefas, incluindo a engenharia de dados<sup>27</sup>. Mas, o campo em que as tecnologias de *Big Data* mais conhecidas são utilizadas é o de processamento de dados para dar suporte à técnicas de mineração de dados e outras atividades de *Data Science*, como pode ser visto na Figura 6 (PROVOST, 2013).

<sup>23</sup> <http://hadoop.apache.org/>

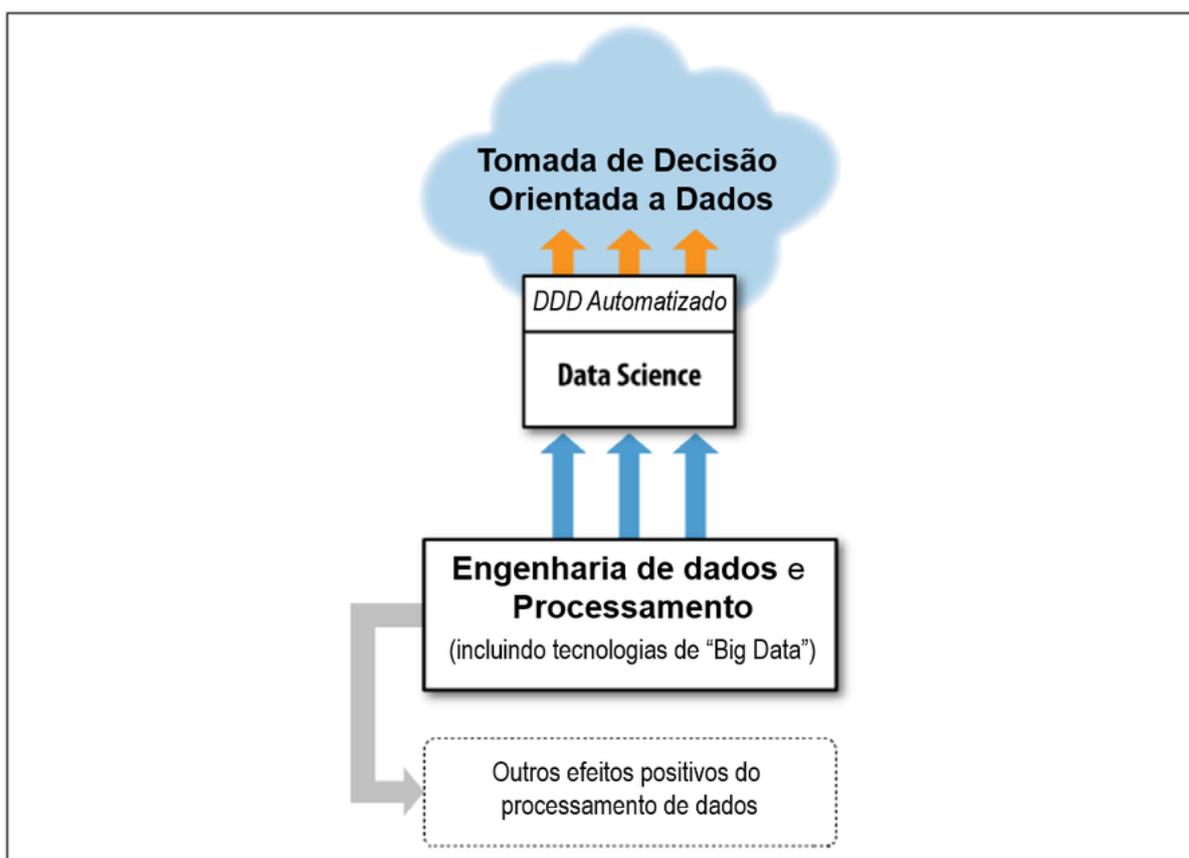
<sup>24</sup> <https://hive.apache.org/>

<sup>25</sup> <https://hbase.apache.org/>

<sup>26</sup> <https://www.mongodb.com/>

<sup>27</sup> área da Engenharia dedicada a processar e tratar dados para aplicações que utilizarão Big Data.

**Figura 6 – Data Science no contexto de vários processos relacionados aos dados de uma organização**



Fonte: Livro "Data Science for Business"

### 2.3.2 Hadoop

Uma plataforma confiável e escalável para armazenamento e análises. É executada em máquinas relativamente baratas e é *open source*<sup>28</sup>, esse é o *Hadoop* (WHITE, 2010). Essa arquitetura proposta como solução para o processamento de *Big Data*, na sua primeira geração, consistia de um sistema distribuído de arquivos – *Hadoop Distributed File System* (HDFS) –, e de um *framework MapReduce*<sup>29</sup> junto com uma interface de coordenação e uma interface para ler e escrever no HDFS (MELOROSE; PERROY; CAREAS, 2015). Essa estrutura fornecia fundamentalmente um sistema de processamento em lotes.

No entanto, desde a sua origem, o *Hadoop* evoluiu além desse tipo de processamento.

O termo "*Hadoop*" algumas vezes é usado para se referir a um grande ecossistema de projetos, não somente o HDFS e o *MapReduce*. Muitos deles são parte

<sup>28</sup> software que possui o seu código fonte aberto e acessível, podendo ser adaptado para diferentes fins.

<sup>29</sup> um framework para processamento distribuído de grandes conjuntos de dados em clusters.

da *Apache Software Foundation*<sup>30</sup>, que provê suporte para a comunidade de projetos *open-source*.

Alguns desses projetos são:

- AVRO<sup>31</sup>;
- Cassandra<sup>32</sup>;
- Chuckwa<sup>33</sup>;
- Hbase<sup>34</sup>;
- Hive<sup>35</sup>;
- Mahout<sup>36</sup>;
- Pig<sup>37</sup>; e
- ZooKeeper<sup>38</sup>.

### 2.3.3 HDFS

Quando um conjunto de dados cresce além da capacidade de armazenamento de uma única máquina, se torna necessário particioná-lo em certo número de máquinas separadas. Sistemas de arquivo que gerenciam o armazenamento por meio de uma rede de computadores são chamados de sistemas de arquivos distribuídos. (WHITE, 2010)

O *Hadoop Distributed File System* (HDFS) é um sistema de arquivos distribuídos, escalável e altamente tolerante a falhas, desenvolvido para ser executado em *hardware* de baixo custo (MELOROSE; PERROY; CAREAS, 2015).

Algumas das características do HDFS são:

- Redundância – o *hardware* pode vir a falhar e os recursos de infra-estrutura podem se esgotar para os processos, mas a redundância presente no HDFS pode lidar com essas situações;

<sup>30</sup> <https://www.apache.org/>

<sup>31</sup> um sistema de serialização de dados.

<sup>32</sup> um banco de dados *multimaster* e escalável.

<sup>33</sup> um sistema de coleção de dados para gerenciar grandes sistemas distribuídos.

<sup>34</sup> um banco de dados distribuído e escalável que suporta o armazenamento de dados estruturados para grandes tabelas.

<sup>35</sup> uma infraestrutura de *data warehouse* que provê a agregação de dados e consultas *ad hoc*.

<sup>36</sup> uma biblioteca de mineração de dados e aprendizado de máquina escalável.

<sup>37</sup> uma linguagem de fluxo de dados de alto nível e um *framework* de execução para computação em paralelo.

<sup>38</sup> um serviço de coordenação de alta performance para aplicações distribuídas.

- Escalabilidade – a escalabilidade linear na camada de armazenamento é necessária para utilizar o processamento paralelo da melhor maneira possível;
- Tolerância a falhas – a habilidade de automaticamente se recuperar da falha e completar o processamento do dado;
- Compatibilidade entre diferentes plataformas – a habilidade de integração entre várias plataformas; e
- Computação e armazenamento em um mesmo ambiente – os dados e a computação são colocados na mesma arquitetura removendo I/O<sup>39</sup> redundantes e acesso excessivo ao disco.

Um *cluster*<sup>40</sup> HDFS possui dois tipos de nós operando no padrão *master-worker* (mestre-trabalhador): um *namenode* (o mestre) e vários *datanodes* (trabalhadores).

O *namenode* gerencia o *namespace*<sup>41</sup> do sistema de arquivos (WHITE, 2010) por meio de operações de abrir, fechar, mover, nomear e renomear arquivos e diretórios (MELOROSE; PERROY; CAREAS, 2015). Além disso, o *namenode* também sabe quais são os *datanodes* que contém cada bloco de um determinado arquivo. Por sua vez, os *datanodes* armazenam e recuperam blocos quando requisitado, e eles, periodicamente, também reportam para o *namenode* uma lista com os blocos que eles estão armazenando.

Sem o *namenode*, o sistema de arquivos não poderia ser mais utilizado. Se a máquina que está executando o *namenode* parasse de funcionar, todos os arquivos presentes no sistema de arquivos seriam perdidos uma vez que não se teria como saber como reconstruir os arquivos que estão particionados pelos *datanodes*. Por esse motivo, o *Hadoop* utiliza de *backups* ou de *namenodes* secundários para evitar esse tipo de problema.

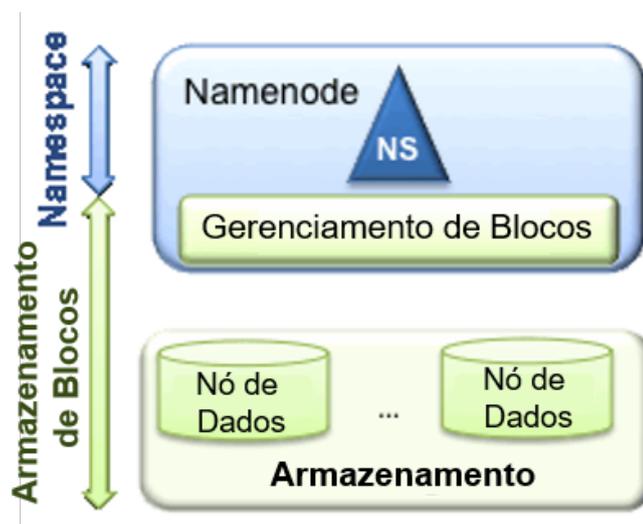
Na Figura 7 podemos observar a estrutura de um *cluster* HDFS com os seus nós divididos nas camadas de *namespace* e *block storage* (armazenamento de bloco).

<sup>39</sup> operações de entrada ou saída de dados por meio de algum código ou programa, para algum outro programa ou hardware.

<sup>40</sup> rede de computadores que trabalham em conjunto para que estes possam ser vistos como um sistema único.

<sup>41</sup> espaço que mantém a árvore do sistema de arquivos e os metadados para todos os arquivos e diretórios pertencentes a árvore.

Figura 7 – Estrutura do HDFS



Fonte: <https://hadoop.apache.org/>

### 2.3.4 MapReduce

MapReduce é um modelo de programação utilizado para processar conjuntos de dados extremamente grandes (MELOROSE; PERROY; CAREAS, 2015). Esse modelo funciona quebrando o processamento em duas fases: a fase de mapeamento e a fase de redução. Cada fase tem pares chave-valor como entrada e saída. Os tipos de cada um deles podem ser definidos pelo programador. O programador também especifica duas funções: a função Map e a função Reduce (WHITE, 2010).

A função *Map* recebe um conjunto de dados e o converte em um outro conjunto de dados, no qual elementos individuais são quebrados em tuplas.

A função *Reduce* utiliza a saída de um *Map* como sua entrada e combina essas tuplas de dados em um conjunto de tuplas menor.

### 2.3.5 Hive

*Hive* é um *framework* para aplicar técnicas de *data warehousing* que realiza o gerenciamento dos dados armazenados no HDFS e provê uma linguagem de consulta baseada em SQL<sup>42</sup>.

O *Hive* surgiu a partir da necessidade de gerenciar e aprender sobre os grandes volumes de dados que o *Facebook*<sup>43</sup> estava produzindo todos os dias na sua rede social. Depois de tentar alguns sistemas diferentes, o time escolheu o *Hadoop* para o

<sup>42</sup> *Structured Query Language* (Linguagem de Consulta Estruturada); linguagem de pesquisa declarativa padrão para bancos de dados relacionais.

<sup>43</sup> <https://www.facebook.com/>

armazenamento e processamento, já que ele tem baixo custo e atende aos requisitos de escalabilidade.

Esse framework foi criado para que analistas com fortes habilidades em SQL pudessem executar consultas em grandes volumes de dados. Hoje, *Hive* é um projeto de sucesso da *Apache* usado por várias empresas como uma plataforma de processamento de dados escalável (WHITE, 2010).

Os principais objetivos do projeto *Hive* são:

- Construir um sistema para gerenciar e consultar dados usando técnicas estruturadas no *Hadoop*;
- Usar *MapReduce* nativo para execução nas camadas do HDFS e do *Hadoop*;
- Usar o HDFS para realizar o armazenamento dos dados do *Hive*;
- Armazenar metadados chave em um RDBMS<sup>44</sup>;
- Estender interfaces SQL, uma ferramenta familiar de *data warehousing* em uso pelas empresa;
- Alta extensibilidade: tipos e funções definidas pelo usuário, formatos, e *scripts*;
- Aumentar a escalabilidade e o desempenho do *Hadoop*; e
- Interoperabilidade com outras plataformas.

Como mencionado anteriormente, essa ferramenta dá suporte a consultas expressadas em uma linguagem declarativa semelhante a SQL – *HiveQL*<sup>45</sup> –, que são compiladas em *MapReduce jobs* executadas no *Hadoop*. Além disso, o *Hive* também inclui um catálogo do sistema, a *metastore*, que contém esquemas e estatísticas, e é usado na exploração de dados e na otimização de consultas.

#### 2.3.5.1 Arquitetura do Hive

A arquitetura do sistema *Hive* pode ser vista na Figura 8. Os principais componentes dessa arquitetura são:

- *Metastore* – armazena o catálogo do sistema e metadados sobre as tabelas, colunas e partições;

<sup>44</sup> Relational Database Management System; é um sistema de gerenciamento de banco de dados que é baseado num modelo relacional.

<sup>45</sup> <https://cwiki.apache.org/confluence/display/Hive/LanguageManual>

- *Driver* – mantém detalhes da sessão, identificadores dos processos, e estatísticas, e gerencia o ciclo de vida de uma declaração *HiveQL* enquanto ela se transporta pelo *Hive*;
- Compilador de consultas – compila *HiveQL* em tarefas de *Map* e *Reduce*;
- Mecanismo de execução – processa e executa as tarefas produzidas pelo compilador em ordem de dependência. O mecanismo de execução gerencia todas as interações entre o compilador e o *Hadoop*;
- *Thrift Server* – provê uma interface *thrift*<sup>46</sup>, um servidor JDBC/ODBC<sup>47</sup>, e uma API rica para integrar *Hive* com outras aplicações;
- CLI e *Web UI* – duas interfaces para o cliente. A interface de linha de comando (CLI) permite execuções em linha de comando e a interface *web* é um console de gerenciamento; e
- Interfaces – interfaces de extensibilidade que incluem as interfaces *SerDe*<sup>48</sup> e *ObjectInspector*<sup>49</sup>, UDF<sup>50</sup>, e UDAF<sup>51</sup> que permitem que o usuário defina suas próprias funções personalizadas.

<sup>46</sup> Camada de comunicação entre processos.

<sup>47</sup> <http://docs.oracle.com/javase/7/docs/technotes/guides/jdbc/bridge.html>

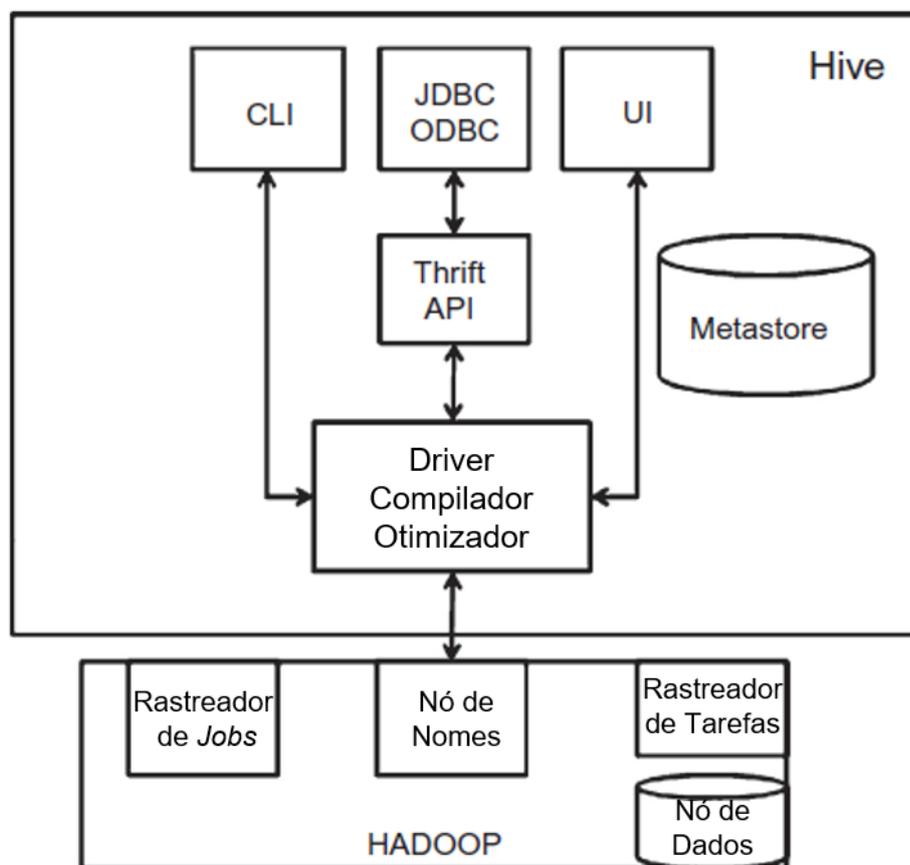
<sup>48</sup> <https://cwiki.apache.org/confluence/display/Hive/SerDe>

<sup>49</sup> <https://hive.apache.org/javadocs/r0.10.0/api/org/apache/hadoop/hive/serde2/objectinspector/ObjectInspector.html>

<sup>50</sup> <https://cwiki.apache.org/confluence/display/Hive/LanguageManual+UDF>

<sup>51</sup> função de agregação definida pelo usuário

Figura 8 – Estrutura do Hive



Fonte: Livro “Data Warehousing in the Age of Big Data”

### 2.3.6 Avro

Apache Avro é um sistema de serialização de dados de linguagem neutra. Esse sistema foi criado com o intuito de resolver a maior desvantagem do *Hadoop*: a falta de portabilidade de linguagens. Ter um formato de dados que pode ser processado por várias linguagens (C<sup>52</sup>, C++<sup>53</sup>, C#<sup>54</sup>, Java<sup>55</sup>, JavaScript<sup>56</sup>, Perl<sup>57</sup>, PHP<sup>58</sup>, Python<sup>59</sup> e Ruby<sup>60</sup>) facilita o compartilhamento dos conjuntos de dados entre usuários (WHITE, 2010).

O Avro provê estruturas de dados ricas; um formato binário de dados rápido e compacto; e chamadas a procedimentos remotos (RPC<sup>61</sup> – *Remote Procedure Calls*)

<sup>52</sup> <https://www.tutorialspoint.com/cprogramming/>

<sup>53</sup> <http://www.cplusplus.com/>

<sup>54</sup> <https://msdn.microsoft.com/en-us/library/67ef8sbd.aspx>

<sup>55</sup> <https://www.oracle.com/java/index.html>

<sup>56</sup> <https://www.javascript.com/>

<sup>57</sup> <https://www.perl.org/>

<sup>58</sup> <http://php.net/>

<sup>59</sup> <https://www.python.org/>

<sup>60</sup> <https://www.ruby-lang.org/pt/>

<sup>61</sup> Tecnologia utilizada para a implementação do modelo cliente-servidor de computação distribuída.

(APACHE, ).

### 2.3.7 Cassandra

Apache Cassandra é um sistema de banco de dados que provê escalabilidade linear<sup>62</sup> e alta disponibilidade sem comprometer o desempenho. Além disso, também é tolerante a falhas em *hardwares* de baixo custo (APACHE, ).

Normalmente, refere-se ao Cassandra como um sistema de banco de dados de arquitetura híbrida, uma vez que este combina um modelo de dados orientado a colunas do *BigTable*<sup>63</sup> com o *MapReduce* do *Hadoop*, e implementa os padrões do *Dynamo*<sup>64</sup> (MELOROSE; PERROY; CAREAS, 2015).

### 2.3.8 Chuckwa

Apache Chuckwa é um sistema de coleção de dados de código aberto utilizado para monitorar grandes sistemas distribuídos. Esse sistema é construído sobre o HDFS e o framework MapReduce, e ele herda a escalabilidade e robustez do Hadoop (APACHE, ).

O Chuckwa também inclui um conjunto de ferramentas flexível utilizado para exibir, monitorar, e analisar resultados para extrair o máximo de informações úteis para o usuário dos dados disponíveis coletados (MELOROSE; PERROY; CAREAS, 2015).

### 2.3.9 HBase

Apache HBase é um sistema de banco de dados distribuído e escalável do Hadoop. Normalmente utilizado quando é necessário realizar acesso randômico de leitura ou escrita em tempo real ao Big Data (APACHE, ).

O HBase também é orientado à colunas e é construído sobre o HDFS (WHITE, 2010). Além disso, esse banco de dados também provê escalabilidade e desempenho ilimitado, e possui certas características de um sistema de banco de dados ACID<sup>65</sup>. O Apache HBase é classificado como um banco de dados *NoSQL* devido a sua arquitetura e o seu *design* serem alinhados com *Base (Being Available and Same Everywhere* – estar disponível e ser o mesmo em qualquer lugar) (MELOROSE; PERROY; CAREAS, 2015).

---

Uma chamada de um procedimento remoto é iniciada pelo cliente enviando uma mensagem para um servidor remoto para que este execute um procedimento específico.

<sup>62</sup> Ao escalonar a aplicação, esta não apresentará falhas nem perdas de recursos, dessa forma o fator de escalabilidade permanecerá constante (THARAKAN, 2007)

<sup>63</sup> <https://cloud.google.com/bigtable/>

<sup>64</sup> [aws.amazon.com/dynamodb](https://aws.amazon.com/dynamodb)

<sup>65</sup> Atomicidade, Consistência, Isolamento e Durabilidade são um conjunto de propriedades das transações em bancos de dados

### 2.3.10 Mahout

Apache Mahout é um projeto cujo objetivo é a construção de um ambiente para a criação rápida de aplicações escaláveis de aprendizado de máquina de elevado desempenho (APACHE, ).

O Mahout consiste de uma coleção algoritmos escaláveis de aprendizado de máquina que são executados no *Hadoop* (SITTO; PRESSER, 2015).

### 2.3.11 Pig

Apache Pig aumenta o nível de abstração para o processamento de grandes conjuntos de dados. Com o Pig, as estruturas de dados são mais ricas, tipicamente sendo multivaloradas e aninhadas, e as transformações que o usuário pode aplicar sobre os dados permitem a extração de informações úteis dos mesmos. (WHITE, 2010).

Em tempo de compilação, a atual infraestrutura do Pig consiste de um compilador que produz sequências de programas *MapReduce*. A arquitetura da linguagem do Pig é uma plataforma de linguagem textual chamada de *Pig Latin*, na qual os objetivos de *design* se basearam no requisito de lidar com grandes processamentos de dados com complexidade mínima (MELOROSE; PERROY; CAREAS, 2015).

### 2.3.12 ZooKeeper

Apache ZooKeeper é uma aplicação utilizada para desenvolver e manter um servidor de código aberto que permita uma coordenação distribuída altamente confiável (APACHE, ).

Em um ambiente controlado como o do RDBMS, as tarefas são geradas de forma controlada e a coordenação precisa somente garantir um gerenciamento da rede sem perda de dados e realizar uma verificação sobre o estado dos nós de um sistema distribuído. No caso do *Hadoop*, o menor volume de dados começa com vários *terabytes* e o dado está distribuído em arquivos pertencentes a diferentes nós. Manter as consultas dos usuários e as tarefas associadas requer um coordenador que seja flexível e escalável, assim como o ZooKeeper (MELOROSE; PERROY; CAREAS, 2015).

## 2.4 Comentários Finais

Os estudos realizados sobre os conceitos de *Big Data*, as tecnologias do *Big Data*, e *Data Warehousing* trouxeram conhecimentos valiosos para uma maior compreensão sobre o estado da arte. Com esse conhecimento, pode-se entender melhor à respeito de aplicações sobre essa área, e possibilita a apresentação de um processo de desenvolvimento de um DW para *Big Data*.

No próximo Capítulo (3) serão apresentados trabalhos relacionados que realizaram a aplicação de técnicas de análise sobre *Big Data* em diversas áreas e cenários. Com isso, será possível observar algumas das possibilidades de desenvolvimento que este tipo de abordagem traz.

### 3 Trabalhos Relacionados

Neste capítulo serão abordados trabalhos relacionados que realizaram a aplicação de técnicas de análise sobre *Big Data*. Para cada trabalho, será apresentado o problema encontrado, a solução aplicada e as limitações, conforme descrito no Quadro 1.

Quadro 1 – Critério de Análise e Apresentação

Atividade	Descrição
Seleção	<ul style="list-style-type: none"> <li>• Selecionar trabalhos que apresentem soluções para problemas relacionados a análises sobre <i>Big Data</i>;</li> <li>• Os trabalhos selecionados precisam definir o problema abordado, as tecnologias aplicadas para a solução, e possíveis limitações ou trabalhos futuros;</li> <li>• Não existe uma filtragem em relação à temática abordada;</li> <li>• Preferência por trabalhos que tenham utilizado algum dos projetos do Apache ou a Cloudera.</li> </ul>
Descrição	<ul style="list-style-type: none"> <li>• Descrever cada um dos trabalhos relacionados apresentando o problema, a solução e as limitações dos mesmos.</li> </ul>

Fonte: Elaborado pelo Autor (2016)

#### 3.1 Trabalho 1 – Análises sobre Imagens Médicas (ZHANG et al., 2016)

##### 3.1.1 Problema

Análises sobre *Big Data* estão chegando à área de saúde. O surgimento das tecnologias da *Internet of Things*<sup>1</sup>, como rastreadores para medição dos sinais vitais, aplicativos de saúde para *smartphones*, estão facilitando o processo de coleta contínua de milhões de usuários sobre os seus dados pessoais de saúde. Uma única instituição pode armazenar cerca de dez a cem *terabytes* de dados sobre imagens médicas.

Além do problema do volume ser claro, também é possível encontrar a necessidade da velocidade pois computações sofisticadas são aplicadas sobre imagens que

<sup>1</sup> Internet das Coisas; é a rede de comunicação entre dispositivos, veículos e outros dispositivos conectados.

surgem a altas taxas de velocidade, e a variedade nos dados também surge quando imagens são posteriormente combinadas com registros estruturados de pacientes.

Análises sobre imagens médicas podem depender de uma *pipeline*<sup>2</sup> altamente especializado para comparar e agregar imagens em 3D de forma única.

Esse trabalho busca entender sobre os comportamentos desse tipo de análise a partir de uma perspectiva de sistemas distribuídos.

### 3.1.2 Solução

A execução do *pipeline* é realizada em paralelo em um *cluster* sobre uma camada de armazenamento distribuída que dá suporte ao compartilhamento de dados entre os nós e entre os estágios do *pipeline*.

Para realizar a execução paralela, foram consideradas duas tecnologias: LSF<sup>3</sup> e *Spark*<sup>4</sup>.

Para a camada de armazenamento foi utilizado o IBM *Spectrum Scale*<sup>5</sup>.

### 3.1.3 Limitações

Após realizar experimentos, foi constatado que realizar tarefas menores faria com que menos memória fosse requisitada, pois a memória foi um dos gargalos encontrados. No entanto, para realizar essa redução, seria necessário aplicar alterações sobre os algoritmos de registro de imagens utilizados e essa tarefa é considerada não trivial.

Além disso, notou-se a necessidade de reescrever o *pipeline* de modo a utilizar o *Spark* de forma mais eficiente. RDDs<sup>6</sup> precisam ser criados para as imagens para que se possa beneficiar do paralelismo oferecido pelo *Spark*.

## 3.2 Trabalho 2 – Transformando o cenário do Mercado com o Cloudera

### 3.2.1 Problema

A EMS (*Experian Marketing Services*) ajuda comerciantes a se conectarem com clientes através da comunicação utilizando uma variedade de canais, guiados por análises avançadas em um extenso banco de dados geográficos, demográficos e de estilo de vida.

<sup>2</sup> Um conjunto de elementos de processamento de dados conectados em série, onde a saída de um elemento é entrada do seguinte.

<sup>3</sup> <http://www-03.ibm.com/systems/spectrum-computing/products/lfs/>

<sup>4</sup> <http://spark.apache.org/>

<sup>5</sup> <http://www-03.ibm.com/systems/storage/spectrum/scale/>

<sup>6</sup> *Resilient Distributed Dataset*; Uma coleção de elementos tolerante a falhas que possam ser operados em paralelo.

A EMS construiu o seu negócio sobre um efetivo processo de coleta, análise, e uso dos dados. No entanto, os comerciantes passaram a exigir mais, pois estes queriam observar o comportamento dos seus respectivos usuários e dar respostas em tempo real. Para isso, seria necessário analisar as últimas compras dos usuários, os padrões de buscas online feitas por eles, e as atividades desses consumidores nas redes sociais.

Mas, a grande quantidade de dados disponíveis nesses canais digitais requer uma infraestrutura tecnológica que possa acomodar processamento rápido, armazenamento em larga escala, e análise flexível sobre dados multiestruturados.

### 3.2.2 Solução

Para definir uma nova infraestrutura foram utilizados dois critérios:

- A infraestrutura deve ter a capacidade de processar tanto dados em tempo real quanto em lote; e
- A infraestrutura deve ter a escalabilidade como uma de suas características para que seja possível acomodar grandes e crescentes volumes de dados.

Tendo esses critérios definidos, a *Cloudera Distribution for Hadoop* (CDH) foi escolhida.

Depois de explorar o *Hadoop*, a EMS criou um mecanismo de acoplamento, chamado de *Cross-Channel Identity Resolution* (CCIR), usado para manter um repositório persistente dos *touchpoints*<sup>7</sup> dos clientes. O CCIR é executado no HBase, que resolve os problemas de persistência, redundância, e tem a habilidade de automaticamente redistribuir os dados. Além disso, o HBase permite tanto o processamento de dados em tempo real quanto em lote.

A Experian realiza o abastecimento de dados no CDH-CCIR usando *scripts* customizados de extração, transformação e carga (ETL) a partir de estruturas internas e de bancos de dados relacionais, incluindo IBM DB2<sup>8</sup>, Oracle<sup>9</sup>, SQL Server<sup>10</sup>. A Experian também usa o *Hive* e o *Pig* para realizar consultas e análises.

### 3.2.3 Limitações

Para a aplicação da nova estrutura envolvendo a CDH foi necessário envolver arquitetos especialistas na área. Além disso, também foi necessário um treinamento

<sup>7</sup> Qualquer forma de interação entre um cliente e um negócio, seja ele pessoalmente, através de uma página *web*, um aplicativo ou qualquer outra forma de comunicação.

<sup>8</sup> <http://www.ibm.com/analytics/us/en/technology/db2/>

<sup>9</sup> <https://www.oracle.com/>

<sup>10</sup> <https://www.microsoft.com/pt-br/server-cloud/products/sql-server/overview.aspx>

sobre *Hadoop* e *HBase* antes de começar a desenvolver a aplicação.

Com a infraestrutura tecnológica que a EMS possuía, a empresa não conseguia acomodar um rápido processamento, um armazenamento em larga escala, e realizar análises flexíveis sobre dados multiestruturados.

Após a implementação do *Hadoop* através da CDH, esse trabalho não apresenta limitações ou trabalhos futuros, além disso este trabalho não explora outras possibilidades tecnológicas que poderiam resolver o problema.

### 3.3 Trabalho 3 – Aprimorando a conectividade de assistências médica com *Big Data*

#### 3.3.1 Problema

Empresas de assistências médicas precisam armazenar os dados sobre saúde por longos períodos. Uma empresa de TI na área de saúde estabeleceu uma política de que se deveria manter um histórico de sete anos dos dados acerca de reclamações e remissões, no entanto os seus sistemas de banco de dados tiveram problemas ao lidar com essa nova política enquanto processava milhares de reclamações todos os dias.

A quantidade de dados era muito grande para que os bancos conseguissem controlá-los. Eles estavam trabalhando acima do limite e estavam sobrecarregados, e isso começou a criar problemas com o processamento de produção em tempo real.

#### 3.3.2 Solução

Utilizando o CDH, a empresa pôde se beneficiar das vantagens oferecidas pelo *Hadoop* que é capaz de separar o *Big Data* dos dados de processamento transacional, e permite um processamento mais suave de informação. Basicamente, ele permite reduzir a sobrecarga presente nos bancos de dados.

Hoje a empresa de TI usa o *Flume*<sup>11</sup> para mover os dados dos seus sistemas para o seu *cluster* CDH todos os dias. A empresa carrega os dados do CDH para um banco de dados da *Oracle* para fins de gerar faturamento, e esse carregamento ocorre uma ou duas vezes por dia utilizando o *Sqoop*<sup>12</sup>.

#### 3.3.3 Limitações

Para poder fazer uso das tecnologias da *Cloudera*, foi necessário avaliar a arquitetura inicial do projeto para que se pudesse tirar maior proveito das capacidades

---

<sup>11</sup> <http://flume.apache.org/>

<sup>12</sup> <http://sqoop.apache.org/>

do *Hadoop*. Além disso, também foi necessário realizar um treinamento sobre as ferramentas para que a solução pudesse ser aplicada.

O trabalho não apresenta a possibilidade de utilização de outras tecnologias que poderiam resolver o problema abordado. Devido a necessidade de uma solução rápida, esse trabalho se limita a algumas tecnologias do ecossistema do *Hadoop*, e não busca apresentar tecnologias alternativas.

### 3.4 Trabalho 4 – Um Sistema Analítico de *Big Data* para analisar a percepção de segurança dos cidadãos (CAMARGO et al., 2016)

#### 3.4.1 Problema

A percepção de segurança (PoS – *Perception of Security*) se refere à avaliação subjetiva dos riscos da população. O PoS é um ponto muito importante para os governantes e os responsáveis pela tomada de decisão ao se planejar a segurança das cidades. Uma prática comum para medir PoS é baseada em pesquisas de opinião, na qual os cidadãos respondem um conjunto de perguntas relacionadas à segurança.

Por um lado, a percepção de segurança depende bastante de muitos aspectos como psicológico, cultural e social, então essa percepção pode mudar dinamicamente quando um desses aspectos muda. Por outro lado, PoS é muito subjetiva já que depende da percepção qualitativa dos indivíduos.

Os governantes precisam ter acesso à informação em tempo real para tomar decisões baseadas em dados atualizados que permitem que os representantes das cidades sejam bem mais precisos. Redes sociais são uma nova forma de comunicação que possuem duas características importantes: informação em tempo real e interatividade.

Nesse trabalho é apresentado um sistema de percepção de segurança que pode ser usado pelos governos a fim de realizar análises sobre dados obtidos a partir de redes sociais com o objetivo de entender, analisar e prevenir eventos de segurança de forma mais sábia.

#### 3.4.2 Solução

A solução é composta por duas camadas. A primeira é a camada de processamento de dados, e a segunda é a camada analítica de dados.

Na primeira camada há cinco componentes:

- *Crawler*<sup>13</sup>;

<sup>13</sup> Programa de computador que navega pela *internet* de uma forma metódica e automatizada.

- Componente de pré-processamento;
- Componente de detecção de segurança;
- Componente de *tagging* automático; e
- Componente de indexação relacionada à segurança.

O *Crawler* permite que os dados sejam coletados a partir do *Twitter*<sup>14</sup>. Uma vez que os dados foram coletados, o componente de pré-processamento extrai campos como usuário, texto do *tweet*, local, e latitude/longitude.

No passo seguinte, os dados são filtrados a fim de detectar somente aqueles que estão relacionados à segurança. O componente de detecção de segurança possui um conjunto de palavras definidas por um especialista na área que são decisivas no processo de filtragem. Uma abordagem de aprendizado de máquina é aplicado de modo a obter um conjunto de dados devidamente classificado.

Usando a biblioteca do *Apache Lucene*<sup>15</sup>, são aplicadas técnicas de indexação baseadas na recuperação de informação que foram construídas como um mecanismo para indexar o significado dos termos relacionados à segurança.

Armazenou-se todos os *tweets* relacionados à segurança em um sistema de banco de dados *NoSQL*<sup>16</sup> chamado de *MongoDB*<sup>17</sup> e em um conjunto de tabelas *Google Fusion*<sup>18</sup>. O *MongoDB* permite uma busca eficiente sobre diferentes campos *JSON*<sup>19</sup>, por esse motivo torna-se possível realizar consultas por data, texto, usuário, e latitude/longitude em cada *tweet*. As tabelas *Google Fusion* permitem a visualização dos *tweets* em um mapa e construção de um mapa de calor.

Com isso, foi possível aplicar análises utilizando mapas e realizando consultas nas bases de modo a obter resultados a respeito da percepção de segurança dos cidadãos.

### 3.4.3 Limitações

A quantidade de dados obtida para realizar o trabalho não era considerável, além disso, o escopo do projeto não englobava um grande número de cidades. O trabalho não discorre sobre possíveis problemas com escalabilidade e também não aponta trabalhos futuros.

<sup>14</sup> <https://twitter.com/>

<sup>15</sup> <https://lucene.apache.org>

<sup>16</sup> <https://aws.amazon.com/pt/nosql/>

<sup>17</sup> <https://www.mongodb.com/>

<sup>18</sup> <https://www.google.com/fusiontables>

<sup>19</sup> <http://www.json.org/>

## 3.5 Trabalho 5 – Aplicando transformações centradas no cliente

### 3.5.1 Problema

A coleta e o uso de dados de forma eficiente se tornou muito importante para a Nokia<sup>20</sup>, uma vez que a empresa busca aprimorar a sua habilidade de entender e prover melhores experiências para os seus usuários com os seus celulares e outros produtos de localização. A companhia utiliza de processamento de dados e análises complexas para construir mapas com modelos de tráfego preditivo e elevação de camadas, para conseguir informações sobre pontos de interesse pelo mundo, para entender sobre a qualidade de seus aparelhos, dentre outras atividades.

Antes da utilização do *Hadoop*, vários grupos da Nokia estavam construindo armazéns para acomodar informações individualmente. Não levou muito tempo para que a empresa percebesse que aquilo poderia proporcionar grande valor caso os seus armazéns viessem a ser integrados, permitindo que os dados globalmente capturados pudessem ser cruzados de modo a obter uma versão da verdade única<sup>21</sup> e compreensível.

### 3.5.2 Solução

Para crescer e dar suporte ao uso extensivo do *Big Data*, a Nokia conta com um ecossistema tecnológico que inclui um *Teradata*<sup>22</sup> *EDW*<sup>23</sup>, vários *data marts* da *Oracle* e do *MySQL*<sup>24</sup>, tecnologias de visualização, e, o principal, o *Hadoop*.

A Nokia possui mais de 100 *terabytes* (TB) de dados estruturados no *Teradata* e *petabytes* (PB) de dados multiestruturados no *HDFS*. O *cluster* centralizado do *Hadoop* que fica no núcleo da infraestrutura da Empresa contém 0.5 PB de dados.

Os *data warehouses* e *data marts* da Nokia continuamente transmitem dados multiestruturados para um ambiente *Hadoop*, permitindo que todos os funcionários da empresa tenham acesso aos dados. A Companhia também faz uso do *Sqoop* para mover os dados do *HDFS* para o *Oracle* ou *Teradata*.

### 3.5.3 Limitações

Inicialmente, a Nokia precisou implementar uma nova infraestrutura que pudesse ter diariamente um fluxo de *terabytes* de dados não estruturados que surgiam a partir

<sup>20</sup> <http://www.nokia.com/>

<sup>21</sup> Do inglês *Single version of truth*, é um conceito técnico utilizado para descrever o *data warehousing* ideal onde todos os dados de uma organização está armazenado de forma consistente e não redundante.

<sup>22</sup> <http://br.teradata.com/?LangType=1046&LangSelect=true>

<sup>23</sup> *Enterprise Data Warehouse*; base de dados unificada que armazena todas as informações de uma empresa e torna as mesmas acessíveis para toda a organização.

<sup>24</sup> <https://www.mysql.com/>

dos aparelhos telefônicos, dos seus serviços, arquivos de *log*, e outras fontes.

Realizar um processamento complexo de dados na escala de *petabytes* usando bancos de dados relacionais era extremamente caro e iria limitar os tipos de dados que poderiam ser consumidos.

Utilizar o *Hadoop* iria facilitar e reduzir bastante os custos. No entanto, para colocar o novo sistema em produção é necessário entender sobre as ferramentas e plataformas que dão suporte ao processo, e estas, por sua vez, tornam a implementação bastante complexa. Por esse motivo, a Nokia passou a utilizar o CDH, que simplifica bastante todo o processo. No entanto, ainda assim, vários meses se passaram até que a empresa decidisse padronizar o uso do *Hadoop*, pois o número de especialistas na ferramenta era muito limitado.

### 3.6 Comentários Finais

Após analisar os cinco trabalhos aqui apresentados, pode-se notar o potencial que a análise sobre *Big Data* possui sobre as mais diversas áreas, e se tornar ciente sobre eventuais problemas relacionados ao desenvolvimento de soluções para esse tipo de análise. O Quadro 2 apresenta um sumário sobre os trabalhos analisados.

Quadro 2 – Sumário de Análises

Trabalho	Tópico	Descrição
Trabalho 1	Problema	<ul style="list-style-type: none"> <li>• Aplicação de análises sobre imagens médicas</li> </ul>
	Tecnologias Utilizadas	<ul style="list-style-type: none"> <li>• LSF</li> <li>• Spark</li> <li>• IBM Spectrum Scale</li> </ul>
	Limitações	<ul style="list-style-type: none"> <li>• Grande requisição sobre memória;</li> <li>• Alterações sobre algoritmos de registro de imagens são muito complexas;</li> <li>• Reescrever <i>pipeline</i> de modo a otimizar o uso do Spark; e</li> <li>• Criar RDDs para se beneficiar do paralelismo oferecido pelo Spark.</li> </ul>
Trabalho 2	Problema	<ul style="list-style-type: none"> <li>• Criação de uma infraestrutura tecnológica que possa acomodar processamento rápido, armazenamento em larga escala, e análise flexível sobre dados multiestruturados</li> </ul>
	Tecnologias Utilizadas	<ul style="list-style-type: none"> <li>• CDH</li> <li>• Hadoop</li> <li>• HBase</li> <li>• IBM DB2</li> <li>• Oracle</li> <li>• SQL Server</li> <li>• Hive</li> <li>• Pig</li> </ul>
	Limitações	<ul style="list-style-type: none"> <li>• Não apresenta limitações ou trabalhos futuros;</li> <li>• O trabalho não descreve outras possibilidades para solucionar o problema.</li> </ul>
Trabalho 3	Problema	<ul style="list-style-type: none"> <li>• Aprimorar a conectividade de assistências médicas</li> </ul>
	Tecnologias Utilizadas	<ul style="list-style-type: none"> <li>• CDH</li> <li>• Hadoop</li> <li>• Flume</li> <li>• Oracle</li> <li>• Sqoop</li> </ul>
	Limitações	<ul style="list-style-type: none"> <li>• Devido a necessidade de uma solução rápida, esse trabalho se limita a algumas tecnologias do ecossistema do <i>Hadoop</i>, e não busca apresentar tecnologias alternativas.</li> </ul>
Trabalho 4	Problema	<ul style="list-style-type: none"> <li>• Analisar a percepção de segurança dos cidadãos</li> </ul>
	Tecnologias Utilizadas	<ul style="list-style-type: none"> <li>• Apache Lucene</li> <li>• MongoDB</li> <li>• Google Fusion</li> </ul>
	Limitações	<ul style="list-style-type: none"> <li>• Quantidade de dados pequena;</li> <li>• O trabalho não discorre sobre possíveis problemas com escalabilidade; e</li> <li>• O trabalho não aponta próximos passos.</li> </ul>
Trabalho 5	Problema	<ul style="list-style-type: none"> <li>• Aprimorar habilidade de entender e prover melhores experiências para os usuários</li> </ul>
	Tecnologias Utilizadas	<ul style="list-style-type: none"> <li>• Teradata EDW</li> <li>• Oracle</li> <li>• MySQL</li> <li>• Hadoop</li> <li>• Sqoop</li> <li>• HDFS</li> </ul>
	Limitações	<ul style="list-style-type: none"> <li>• Utilização do Hadoop não ocorreu de imediato devido à complexidade de uso da ferramenta;</li> <li>• Foi necessário a utilização da CDH para implantar a nova infraestrutura;</li> <li>• Por não ter muitos especialistas na área, o tempo de implantação tornou-se mais longo.</li> </ul>

Além disso, também se pode observar quais tecnologias foram utilizadas e entender o porque destas terem sido selecionadas.

No próximo capítulo será apresentado um processo para o desenvolvimento de um *Data Warehouse* para um cenário de *Big Data* utilizando a ferramenta *Hadoop*.

## 4 Arquitetura e Processo de Desenvolvimento do DW

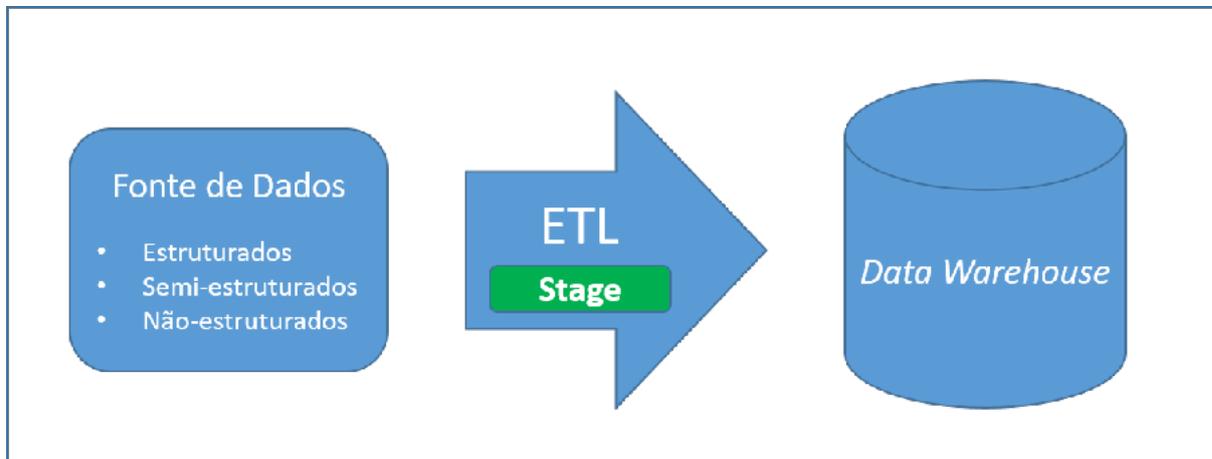
Neste capítulo serão apresentados a arquitetura e o processo de desenvolvimento do DW. Além disso, também será descrita a especificação das ferramentas utilizadas para a construção do *data warehouse*.

### 4.1 Arquitetura do DW

A arquitetura do DW (Figura 9) desenvolvida neste trabalho consiste dos seguintes elementos:

- Fonte de Dados
- ETL
- *Data Warehouse*

Figura 9 – Arquitetura do Data Warehousing



Fonte: Elaborada pelo Autor (2016)

A fonte de dados abrange todas as informações que serão obtidas a partir da origem. Estes dados são classificados em estruturados, semiestruturados e não-estruturados. Durante o processo de ETL, os dados então são classificados e categorizados de modo a dar significado para os mesmos. Por fim, com as devidas transformações aplicadas e com os dados limpos, é possível carregá-los no DW.

## 4.2 Especificação das Ferramentas

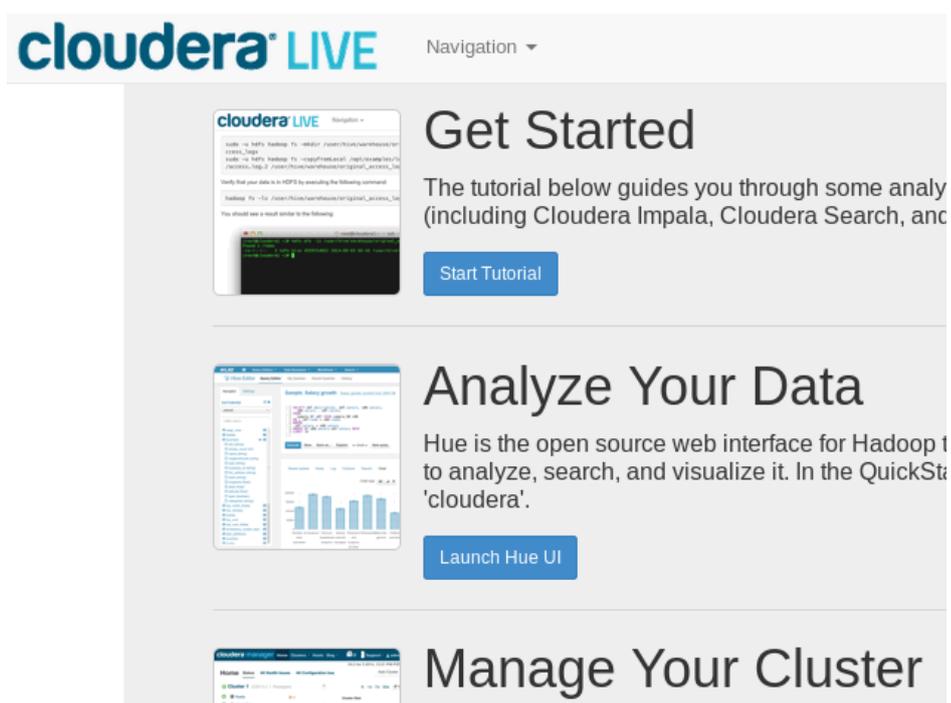
Para a construção do DW foi utilizada a plataforma *Hadoop*, o *framework Hive*, a máquina virtual da *Cloudera Distribution for Hadoop*, e o *framework* de geração de dados da TPCx-BB.

Por ser uma plataforma de código aberto e ser relativamente barata, o *Hadoop* foi selecionado para apresentar o processamento de *Big Data*.

A fim de desenvolver um DW, O *Hive* foi selecionado, pois este *framework* foi criado com o intuito de aplicar técnicas de *data warehousing* para realizar o gerenciamento dos dados armazenados no HDFS. Além disso, o *HiveQL* – linguagem de consulta do *Hive* – possui muitas semelhanças com SQL, o que facilita e auxilia bastante no desenvolvimento e aprendizado.

A CDH foi utilizada pelo fato de ser uma plataforma de distribuição completamente *open source* que integra o *Hadoop* com vários outros projetos incluindo o *Hive*. A máquina virtual da Cloudera já traz todas as tecnologias instaladas e prontas para uso, além de facilitar todo o processo de manutenção e gerenciamento através de suas interfaces (Figura 10). O CDH foi utilizado no modo pseudo distribuído, pois este é recomendado para aprendizado sobre as tecnologias envolvidas pelo fato de não ser necessário utilizar um *cluster* de máquinas e conseguir simular um cluster com apenas uma máquina.

Figura 10 – Interfaces da Cloudera



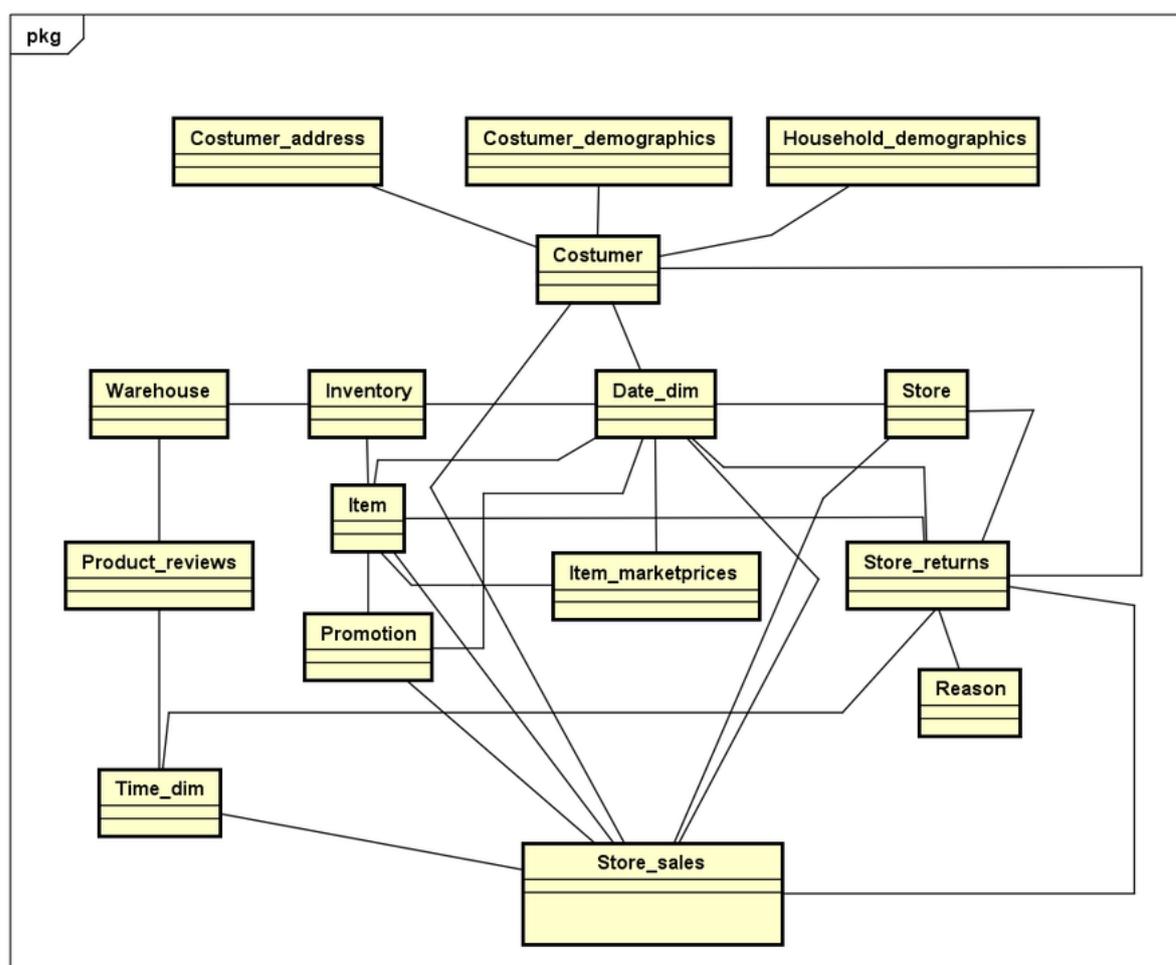
Fonte: Elaborada pelo Autor (2016)

Por fim, o TPCx-BB foi utilizado, pois este define referências de bancos de dados e processamento de transações *Big Data*. O seu *framework* de geração de dados possui parâmetros que permitem o escalonamento de dados e torna possível aplicar testes em diversos cenários. Desse modo, ao utilizar a CDH no modo pseudo distribuído pode-se também definir um volume de dados que seja suficiente para aplicar os testes e compreender sobre o funcionamento do processo sem haver necessidade de se preocupar com limitações de *hardware*.

### 4.2.1 Desenvolvimento do DW

Para desenvolver o DW, este trabalho utilizou a estrutura das tabelas geradas pelo TPCx-BB. Na Figura 11 pode-se observar um diagrama de classes e ver como as tabelas se relacionam.

Figura 11 – Diagrama de Classes



powered by Astah

Fonte: Elaborada pelo Autor (2016)

A seguir, serão descritas as tabelas que compõem o *data warehouse*, o qual foi

construído sobre um modelo de varejo composto por 15 dimensões e 3 fatos.

#### 4.2.1.1 Dimensões

##### 4.2.1.1.1 *Data\_dim*

**Tabela 1 – Dimensão de Datas.**

Nome	Tipo	Descrição
<b>d_date_sk</b>	BIGINT	Chave começando em 1 <i>String</i> única de tamanho 16
<b>d_date_id</b>	CHAR(16)	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>d_date</b>	CHAR(13)	aaaa-MM-dd
<b>d_month_seq</b>	INT	Começa em 0 Conta todos os meses desde a primeira até a última data existente na tabela Começa em 0
<b>d_week_seq</b>	INT	Começa em 0 Conta todas as semanas desde a primeira até a última data existente na tabela Começa em 0
<b>d_quarter_seq</b>	INT	Começa em 0 Conta todos os trimestres desde a primeira até a última data existente na tabela
<b>d_year</b>	INT	Parte da data referente ao ano: aaaa
<b>d_dow</b>	INT	Dia da semana: 1-7, 1 = Segunda
<b>d_moy</b>	INT	Mês do ano: 1-12, 1 = Janeiro
<b>d_dom</b>	INT	Dia do mês: 1-31
<b>d_qoy</b>	INT	Trimestre do ano: 1-4
<b>d_fy_year</b>	INT	Ano financeiro: $d\_year + 1/2year$
<b>d_fy_quarter_seq</b>	INT	Trimestre financeiro: $d\_quarter + 1/2year$
<b>d_fy_week_seq</b>	INT	Semana financeira: $d\_week + 1/2year$
<b>d_day_name</b>	CHAR(9)	Dia da semana como uma <i>string</i> : {Segunda, ..., Domingo}

Nome	Tipo	Descrição
<b>d_quarter_name</b>	CHAR(6)	Trimestre do ano como uma <i>string</i> : <i>aaaaT{1..4}</i> : exemplo: 1990T2
<b>d_holiday</b>	CHAR(1)	Feriado (N/Y)
<b>d_weekend</b>	CHAR(1)	Fim de Semana (N/Y)
<b>d_following_holiday</b>	CHAR(1)	Feriado Seguinte (N/Y)
<b>d_first_dom</b>	INT	Primeiro dia do mês no calendário Juliano
<b>d_last_dom</b>	INT	Último dia do mês no calendário Juliano
<b>d_same_day_ly</b>	INT	O dia corrente no calendário Juliano
<b>d_same_day_lq</b>	INT	O dia corrente no calendário Juliano
<b>d_current_day</b>	CHAR(1)	Dia corrente (N/Y)
<b>d_current_week</b>	CHAR(1)	Semana corrente (N/Y)
<b>d_current_month</b>	CHAR(1)	Mês corrente (N/Y)
<b>d_current_quarter</b>	CHAR(1)	Trimestre corrente (N/Y)
<b>d_current_year</b>	CHAR(1)	Ano corrente (N/Y)

#### 4.2.1.1.2 Cliente

**Tabela 2 – Dimensão de Clientes.**

Nome	Tipo	Descrição
<b>c_customer_sk</b>	<b>BIGINT</b>	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>c_customer_id</b>	<b>CHAR(16)</b>	Conjunto de caracteres: "ABCDEFGHIJKLMNOPQRSTUVWXYZ"
<b>c_current_cdemo_sk</b>	<b>BIGINT</b>	Referência aleatória para a tabela <i>customer_demographics cd_demo_sk</i>
<b>c_current_hdemo_sk</b>	<b>BIGINT</b>	Referência aleatória para a tabela <i>household_demographics hd_demo_sk</i>
<b>c_current_addr_sk</b>	<b>BIGINT</b>	Referência aleatória para a tabela <i>customer_address ca_address_sk</i>
<b>c_first_shipto_date_sk</b>	<b>BIGINT</b>	Referência aleatória para a tabela <i>date_dim</i> <i>d_date_sk</i>
<b>c_first_sales_date_sk</b>	<b>BIGINT</b>	Referência aleatória para a tabela <i>date_dim</i> <i>d_date_sk</i>

Nome	Tipo	Descrição
<b>c_salutation</b>	<b>CHAR(10)</b>	Saudação {Mr., Mrs., Ms., Miss., Sir., Dr.}
<b>c_first_name</b>	<b>CHAR(20)</b>	Primeiro nome
<b>c_last_name</b>	<b>CHAR(30)</b>	Último nome
<b>c_preferred_cust_flag</b>	<b>CHAR(1)</b>	<i>Flag</i> para determinar se esse é um cliente “preferencial”
<b>c_birth_day</b>	<b>INT</b>	Número aleatório: [1, 31]
<b>c_birth_month</b>	<b>INT</b>	Número aleatório: [1, 12]
<b>c_birth_year</b>	<b>INT</b>	Número aleatório: [1924, 1992]
<b>c_birth_country</b>	<b>VARCHAR(20)</b>	País natal
<b>c_login</b>	<b>CHAR(13)</b>	<i>String</i> aleatória, tamanho: [1-13] Padrão:
<b>c_email_address</b>	<b>CHAR(50)</b>	C_first_name.c_last_name@provedorRandômico.tld Data do ultimo <i>review</i>
<b>c_last_review_date</b>	<b>VARCHAR(19)</b>	Min: Dia atual – 1 ano Max: Dia atual

#### 4.2.1.1.3 Endereco\_cliente

**Tabela 3 – Endereço do cliente**

Nome	Tipo	Descrição
<b>ca_address_sk</b>	<b>BIGINT</b>	<i>Serial key</i> String única, tamanho: 16
<b>ca_address_id</b>	<b>CHAR(16)</b>	Conjunto de caracteres: "ABCDEFGHIJKLMNOPQRSTUVWXYZ" Nº da rua
<b>ca_street_number</b>	<b>CHAR(10)</b>	Valor randômico: [1, 1000]
<b>ca_street_name</b>	<b>VARCHAR(60)</b>	Nome da rua Tipo da rua
<b>ca_street_type</b>	<b>CHAR(15)</b>	{Street,ST,Avenue,Ave,Boulevard, Blvd,Road,RD,Parkway,Pkwy,Way, Wy,Drive,Dr.,Circle,Cir.,Lane,Ln ,Court,Ct.}
<b>ca_suite_number</b>	<b>CHAR(10)</b>	String aleatória, tamanho: [1, 10]
<b>ca_city</b>	<b>VARCHAR(60)</b>	Cidades
<b>ca_county</b>	<b>VARCHAR(30)</b>	Condado
<b>ca_state</b>	<b>CHAR(2)</b>	Estado
<b>ca_zip</b>	<b>CHAR(10)</b>	Código postal
<b>ca_country</b>	<b>VARCHAR(20)</b>	País
<b>ca_gmt_offset</b>	<b>DECIMAL</b>	Fuso horário
<b>ca_location_type</b>	<b>CHAR(20)</b>	Tipo de Lugar {Single family, Condo, Apartment}

#### 4.2.1.1.4 Demografia\_cliente

Tabela 4 – Dados demográficos do cliente

Nome	Tipo	Descrição
<b>cd_demo_sk</b>	<b>BIGINT</b>	<i>Serial key</i>
<b>cd_gender</b>	<b>CHAR(1)</b>	Masculino ou Feminino (M/F)
<b>cd_marital_status</b>	<b>CHAR(1)</b>	Estado Civil
<b>cd_education_status</b>	<b>CHAR(20)</b>	Escolaridade

Nome	Tipo	Descrição
<b>cd_purchase_estimate</b>	<b>INT</b>	Poder de compra
<b>cd_credit_rating</b>	<b>CHAR(10)</b>	Avaliação de crédito
<b>cd_dep_count</b>	<b>INT</b>	Dependentes
<b>cd_dep_employed_count</b>	<b>INT</b>	Dependentes Empregados
<b>cd_dep_college_count</b>	<b>INT</b>	Dependentes na faculdade

#### 4.2.1.1.5 Demografia\_domicilio

**Tabela 5 – Dados demográficos do domicílio**

Nome	Tipo	Descrição
<b>hd_demo_sk</b>	<b>BIGINT</b>	<i>Serial key</i>
<b>hd_buy_potential</b>	<b>CHAR(15)</b>	Poder de compra
<b>hd_dep_count</b>	<b>INT</b>	Dependentes
<b>hd_vehicle_count</b>	<b>INT</b>	Quantidade de Veículos

#### 4.2.1.1.6 Inventario

**Tabela 6 – Inventário dos itens estocados no warehouse**

Nome	Tipo	Descrição
<b>inv_date_sk</b>	<b>BIGINT</b>	Referência randômica para tabela <i>date_dimd_date_sk</i>
<b>inv_item_sk</b>	<b>BIGINT</b>	Referência randômica para tabela <i>item i_item_sk</i>
<b>inv_warehouse_sk</b>	<b>BIGINT</b>	Referência randômica para tabela <i>warehouse w_warehouse_sk</i>
<b>inv_quantity_on_hand</b>	<b>INT</b>	Quantidade no inventário

#### 4.2.1.1.7 Item

Tabela 7 – Dimensão Item

Nome	Tipo	Descrição
<b>i_item_sk</b>	BIGINT	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>i_item_id</b>	CHAR(16)	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>i_rec_start_date</b>	CHAR(13)	Data de início de disponibilidade
<b>i_rec_end_date</b>	CHAR(13)	Data de fim de disponibilidade
<b>i_item_desc</b>	VARCHAR(200)	Descrição do item Preço atual
<b>i_current_price</b>	DECIMAL	Número decimal randômico: [0.09, 99.99]
<b>i_wholesale_cost</b>	DECIMAL	Custo total de venda
<b>i_brand_id</b>	INT	Código da marca
<b>i_brand</b>	CHAR(50)	Marca
<b>i_class_id</b>	INT	Código da classe
<b>i_class</b>	CHAR(50)	Classe
<b>i_category_id</b>	INT	Código da categoria
<b>i_category</b>	CHAR(50)	Categoria
<b>i_manufact_id</b>	INT	Código da manufatura
<b>i_manufact</b>	CHAR(50)	Manufatura
<b>i_size</b>	CHAR(20)	Tamanho do item
<b>i_formulation</b>	CHAR(20)	Formulação
<b>i_color</b>	CHAR(20)	Cor
<b>i_units</b>	CHAR(10)	Unidades
<b>i_container</b>	CHAR(10)	Recipiente
<b>i_manager_id</b>	INT	Código do gerente
<b>i_product_name</b>	CHAR(50)	Nome do produto

## 4.2.1.1.8 Revisoes\_Produtos

Tabela 8 – Dimensão com informações sobre revisões dos produtos

Nome	Tipo	Descrição
------	------	-----------

Nome	Tipo	Descrição
<b>pr_review_sk</b>	BIGINT	<i>Serial key</i>
<b>pr_review_date</b>	CHAR(13)	Data do <i>review</i>
<b>pr_review_time</b>	CHAR(6)	Referência randômica para a tabela <i>time_dim</i> <i>t_time_sk</i>
<b>pr_review_rating</b>	INT	Nota
<b>pr_item_sk</b>	BIGINT	Referência randômica para <i>ws_item_sk</i>
<b>pr_user_sk</b>	BIGINT	Referência randômica para <i>ws_user_sk</i>
<b>pr_order_sk</b>	BIGINT	Referência randômica para a tabela <i>web_sales</i> <i>order_id</i>
<b>pr_review_content</b>	STRING	Descrição da <i>review</i>

4.2.1.1.9 *Promocao*

Tabela 9 – Dimensão com informações sobre as promoções

Nome	Tipo	Descrição
<b>p_promo_sk</b>	BIGINT	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>p_promo_id</b>	CHAR(16)	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>p_start_date_sk</b>	BIGINT	Referência randômica para a tabela <i>date_dim</i> <i>d_date_sk</i>
<b>p_end_date_sk</b>	BIGINT	Referência randômica para a tabela <i>date_dim</i> <i>d_date_sk</i>
<b>p_item_sk</b>	BIGINT	Referência randômica para a tabela <i>item</i> <i>i_item_sk</i>
<b>p_cost</b>	DECIMAL	Custo
<b>p_promo_name</b>	CHAR(50)	Nome da promoção
<b>p_channel_dmail</b>	CHAR(1)	Y/N
<b>p_channel_email</b>	CHAR(1)	Y/N

Nome	Tipo	Descrição
<b>p_channel_catalog</b>	CHAR(1)	Y/N
<b>p_channel_tv</b>	CHAR(1)	Y/N
<b>p_channel_radio</b>	CHAR(1)	Y/N
<b>p_channel_press</b>	CHAR(1)	Y/N
<b>p_channel_event</b>	CHAR(1)	Y/N
<b>p_channel_demo</b>	CHAR(1)	Y/N
<b>p_channel_details</b>	VARCHAR(100)	Detalhes sobre o canal
<b>p_purpose</b>	CHAR(15)	Propósito
<b>p_discount_active</b>	CHAR(1)	Y/N

#### 4.2.1.1.10 Razao

**Tabela 10 – Dimensão com informações sobre as razões de retorno dos itens**

Nome	Tipo	Descrição
<b>r_reason_sk</b>	BIGINT	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>r_reason_id</b>	CHAR(16)	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>r_reason_desc</b>	CHAR(100)	Descrição da razão do retorno

#### 4.2.1.1.11 Loja

**Tabela 11 – Dimensão de lojas**

Nome	Tipo	Descrição
<b>s_store_sk</b>	BIGINT	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>s_store_id</b>	CHAR(16)	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”

Nome	Tipo	Descrição
<b>s_rec_start_date</b>	CHAR(13)	Referência randômica para a tabela <i>date_dim d_date_sk</i>
<b>s_rec_end_date</b>	CHAR(13)	Referência randômica para a tabela <i>date_dim d_date_sk</i>
<b>s_closed_date_sk</b>	BIGINT	Data em que a loja está fechada. Referência para a tabela <i>date_dim d_date_sk</i>
<b>s_store_name</b>	VARCHAR(50)	Nome da loja
<b>s_number_employees</b>	INT	Número inteiro randômico: [200, 300]
<b>s_floor_space</b>	INT	Número inteiro randômico: [5000000, 10000000]
<b>s_hours</b>	CHAR(20)	Horário de atendimento
<b>s_manager</b>	VARCHAR(40)	Nome do Gerente
<b>s_market_id</b>	INT	Código do Mercado
<b>s_geography_class</b>	VARCHAR(100)	Classe geográfica
<b>s_market_desc</b>	VARCHAR(100)	Descrição do Mercado
<b>s_market_manager</b>	VARCHAR(40)	Gerente do Mercado
<b>s_division_id</b>	INT	Código da Divisão
<b>s_division_name</b>	VARCHAR(50)	Divisão
<b>s_company_id</b>	INT	Código da Empresa
<b>s_company_name</b>	VARCHAR(50)	Empresa
<b>s_street_number</b>	VARCHAR(10)	Endereço igual ao do <i>warehouse</i>
<b>s_street_name</b>	VARCHAR(60)	Endereço igual ao do <i>warehouse</i>
<b>s_street_type</b>	CHAR(15)	Endereço igual ao do <i>warehouse</i>
<b>s_suite_number</b>	CHAR(10)	Endereço igual ao do <i>warehouse</i>
<b>s_city</b>	VARCHAR(60)	Endereço igual ao do <i>warehouse</i>
<b>s_county</b>	VARCHAR(30)	Endereço igual ao do <i>warehouse</i>
<b>s_state</b>	CHAR(2)	Endereço igual ao do <i>warehouse</i>
<b>s_zip</b>	CHAR(10)	Endereço igual ao do <i>warehouse</i>
<b>s_country</b>	VARCHAR(20)	Endereço igual ao do <i>warehouse</i>
<b>s_gmt_offset</b>	DECIMAL	Endereço igual ao do <i>warehouse</i>
<b>s_tax_percentage</b>	DECIMAL	Taxa

#### 4.2.1.1.12 *Tempo\_dim*

Tabela 12 – Dimensão de tempo

Nome	Tipo	Descrição
<b>t_time_sk</b>	<b>BIGINT</b>	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>t_time_id</b>	<b>CHAR(16)</b>	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ” Começa no 0
<b>t_time</b>	<b>INT</b>	$t\_time = t\_time\_id$
<b>t_hour</b>	<b>INT</b>	$t\_time\_id/60/60 \bmod 24$
<b>t_minute</b>	<b>INT</b>	$t\_time\_id/60 \bmod 60$
<b>t_second</b>	<b>INT</b>	$t\_time\_id \bmod 60$
<b>t_am_pm</b>	<b>CHAR(2)</b>	Turno {AM,PM}
<b>t_shift</b>	<b>CHAR(20)</b>	{ <i>first, second, third</i> }
<b>t_sub_shift</b>	<b>CHAR(20)</b>	{ <i>morning, afternoon, night</i> }
<b>t_meal_time</b>	<b>CHAR(20)</b>	{ <i>breakfast, lunch</i> }

## 4.2.1.1.13 Armazem

Tabela 13 – Dimensão de armazéns

Nome	Tipo	descrição
<b>w_warehouse_sk</b>	<b>BIGINT</b>	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>w_warehouse_id</b>	<b>CHAR(16)</b>	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>w_warehouse_name</b>	<b>VARCHAR(20)</b>	Nome do <i>warehouse</i>
<b>w_warehouse_sq_ft</b>	<b>INT</b>	
<b>w_street_number</b>	<b>CHAR(10)</b>	Número da rua
<b>w_street_name</b>	<b>VARCHAR(60)</b>	Nome da rua
<b>w_street_type</b>	<b>CHAR(15)</b>	Tipo da rua
<b>w_suite_number</b>	<b>CHAR(10)</b>	
<b>w_city</b>	<b>VARCHAR(60)</b>	Cidade

Nome	Tipo	descrição
<b>w_county</b>	<b>VARCHAR(30)</b>	Condado
<b>w_state</b>	<b>CHAR(2)</b>	Estado
<b>w_zip</b>	<b>CHAR(10)</b>	Código Postal
<b>w_country</b>	<b>VARCHAR(20)</b>	País
<b>w_gmt_offset</b>	<b>DECIMAL</b>	Fuso Horário

#### 4.2.1.1.14 Pagina\_Web

**Tabela 14 – Dimensão de Página da Web**

Nome	Tipo	Descrição
<b>wp_web_page_sk</b>	<b>BIGINT</b>	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>wp_web_page_id</b>	<b>CHAR(16)</b>	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>wp_rec_start_date</b>	<b>CHAR(13)</b>	Referência para tabela <i>date_dim d_date_sk</i>
<b>wp_rec_end_date</b>	<b>CHAR(13)</b>	Referência para tabela <i>date_dim d_date_sk</i>
<b>wp_creation_date_sk</b>	<b>BIGINT</b>	Referência para tabela <i>date_dim d_date_sk</i>
<b>wp_access_date_sk</b>	<b>BIGINT</b>	Referência para tabela <i>date_dim d_date_sk</i>
<b>wp_autogen_flag</b>	<b>CHAR(1)</b>	
<b>wp_customer_sk</b>	<b>BIGINT</b>	Referência para tabela <i>customer c_customer_sk</i>
<b>wp_url</b>	<b>VARCHAR(100)</b>	URL
<b>wp_type</b>	<b>CHAR(50)</b>	Tipo da página {general,order,welcome,ad,feedback,protected,dynamic,review}
<b>wp_char_count</b>	<b>INT</b>	Quantidade de caracteres na página
<b>wp_link_count</b>	<b>INT</b>	Quantidade de <i>links</i> na página
<b>wp_image_count</b>	<b>INT</b>	Quantidade de imagens na página
<b>wp_max_ad_count</b>	<b>INT</b>	Quantidade máxima de propagandas na página

#### 4.2.1.1.15 Web\_site

Tabela 15 – Dimensão de Web Site

Nome	Tipo	Descrição
<b>web_site_sk</b>	<b>BIGINT</b>	<i>Serial key</i> <i>String</i> única, tamanho: 16
<b>web_site_id</b>	<b>CHAR(16)</b>	Conjunto de caracteres: “ABCDEFGHIJKLMNOPQRSTUVWXYZ”
<b>web_rec_start_date</b>	<b>CHAR(13)</b>	Referência para a tabela <i>date_dim d_date_sk</i>
<b>web_rec_end_date</b>	<b>CHAR(13)</b>	Referência para a tabela <i>date_dim d_date_sk</i>
<b>web_name</b>	<b>VARCHAR(50)</b>	Nome do site
<b>web_open_date_sk</b>	<b>BIGINT</b>	Referência para a tabela <i>date_dim d_date_sk</i>
<b>web_close_date_sk</b>	<b>BIGINT</b>	Referência para a tabela <i>date_dim d_date_sk</i>
<b>web_class</b>	<b>VARCHAR(50)</b>	Classe do site
<b>web_manager</b>	<b>VARCHAR(40)</b>	Gerente do site
<b>web_mkt_id</b>	<b>INT</b>	Código de Mercado do site
<b>web_mkt_class</b>	<b>VARCHAR(50)</b>	Classe do Mercado do site
<b>web_mkt_desc</b>	<b>VARCHAR(100)</b>	Descrição do Mercado do site
<b>web_market_manager</b>	<b>VARCHAR(40)</b>	Gerente do Mercado do site
<b>web_company_id</b>	<b>INT</b>	Código da Empresa do site
<b>web_company_name</b>	<b>CHAR(50)</b>	Empresa do site
<b>web_street_number</b>	<b>CHAR(10)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_street_name</b>	<b>VARCHAR(60)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_street_type</b>	<b>CHAR(15)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_suite_number</b>	<b>CHAR(10)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_city</b>	<b>VARCHAR(60)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_county</b>	<b>VARCHAR(30)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_state</b>	<b>CHAR(2)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_zip</b>	<b>CHAR(10)</b>	Endereço igual ao do <i>warehouse</i>
<b>web_country</b>	<b>VARCHAR(20)</b>	Endereço igual ao do <i>warehouse</i>

Nome	Tipo	Descrição
<b>web_gmt_offset</b>	<b>DECIMAL</b>	Endereço igual ao do <i>warehouse</i>
<b>web_tax_percentage</b>	<b>DECIMAL</b>	Taxa

#### 4.2.1.2 Fatos

##### 4.2.1.2.1 *Item\_precos\_mercado*

**Tabela 16 – Fato Preços de Mercado dos Itens**

Nome	Tipo	Descrição
<b>imp_sk</b>	BIGINT	<i>Serial key</i>
<b>imp_item_sk</b>	BIGINT	Referência randômica para a tabela <i>item</i>
<b>imp_competitor</b>	VARCHAR(20)	<i>i_item_sk</i>
<b>imp_competitor_price</b>	DECIMAL	<i>String</i> randômica, tamanho: [1, 20]
<b>imp_start_date_sk</b>	BIGINT	Decimal randômico [0.99, 99.99]
<b>imp_end_date_sk</b>	BIGINT	Referência randômica para a tabela <i>date_dim</i>
		<i>d_date_sk</i>
		Referência randômica para a tabela <i>date_dim</i>
		<i>d_date_sk</i>

##### 4.2.1.2.2 Retorno\_Vendas

**Tabela 17 – Fato Retornos de Vendas**

Nome	Tipo	descrição
<b>sr_returned_date_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>date_dim</i>
		<i>d_date_sk</i>

Nome	Tipo	descrição
<b>sr_return_time_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>time_dim</i> <i>t_time_sk</i>
<b>sr_item_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>item_i_item_sk</i>
<b>sr_customer_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>costumer</i> <i>c_costumer_sk</i>
<b>sr_cdemo_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>costumer_demographics</i> <i>cd_demo_sk</i>
<b>sr_hdemo_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>household_demographics</i> <i>hd_hdemo_sk</i>
<b>sr_addr_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>costumer_address</i> <i>ca_addr_sk</i>
<b>sr_store_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>store s_store_sk</i>
<b>sr_reason_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>reason</i> <i>r_reason_sk</i>
<b>sr_ticket_number</b>	<b>BIGINT</b>	Referência para a tabela <i>store_sales</i> <i>ss_ticket_number</i>
<b>sr_return_quantity</b>	<b>INT</b>	Quantidade de itens retornados
<b>sr_return_amt</b>	<b>DECIMAL</b>	Preço * Quantidade
<b>sr_return_tax</b>	<b>DECIMAL</b>	<i>sr_return_amt</i> * Taxa
<b>sr_return_amt_inc_tax</b>	<b>DECIMAL</b>	<i>sr_return_amt</i> + <i>sr_return_tax</i>
<b>sr_fee</b>	<b>DECIMAL</b>	
<b>sr_return_ship_cost</b>	<b>DECIMAL</b>	
<b>sr_refunded_cash</b>	<b>DECIMAL</b>	Dinheiro retornado
<b>sr_reversed_charge</b>	<b>DECIMAL</b>	
<b>sr_store_credit</b>	<b>DECIMAL</b>	

Nome	Tipo	descrição
<b>sr_net_loss</b>	<b>DECIMAL</b>	

## 4.2.1.2.3 Vendas\_Loja

Tabela 18 – Fato Vendas da Loja

Nome	Tipo	descrição
<b>ss_sold_date_sk</b>	<b>BIGINT</b>	Referência para a tabela <i>date_dim d_date_sk</i>
<b>ss_sold_time_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>time_dim t_time_sk</i>
<b>ss_item_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>item i_item_sk</i>
<b>ss_customer_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>customer c_customer_sk</i>
<b>ss_cdemo_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>customer_demographics cd_cdemo_sk</i>
<b>ss_hdemo_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>household_demographics hd_hdemo_sk</i>
<b>ss_addr_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>customer_address ca_addr_sk</i>
<b>ss_store_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>store s_store_sk</i>
<b>ss_promo_sk</b>	<b>BIGINT</b>	Referência randômica para a tabela <i>promotion p_promo_sk</i>
<b>ss_ticket_number</b>	<b>BIGINT</b>	Número do <i>ticket</i>
<b>ss_quantity</b>	<b>INT</b>	Quantidade comprada do item

Nome	Tipo	descrição
<b>ss_warehouse_cost</b>	<b>DECIMAL</b>	Custo de venda total
<b>ss_list_price</b>	<b>DECIMAL</b>	Preço de lista de um único item
<b>ss_sales_price</b>	<b>DECIMAL</b>	Preço de venda de um único item
<b>ss_ext_discount_amt</b>	<b>DECIMAL</b>	Desconto * Quantidade
<b>ss_ext_sales_price</b>	<b>DECIMAL</b>	Custo do item * Quantidade
<b>ss_ext_warehouse_cost</b>	<b>DECIMAL</b>	Custo de venda total * Quantidade
<b>ss_ext_list_price</b>	<b>DECIMAL</b>	Preço de lista do item * Quantidade
<b>ss_ext_tax</b>	<b>DECIMAL</b>	Taxa
<b>ss_coupon_amt</b>	<b>DECIMAL</b>	Cupom de desconto
<b>ss_net_paid</b>	<b>DECIMAL</b>	Valor líquido pago dos itens * Quantidade
<b>ss_net_paid_inc_tax</b>	<b>DECIMAL</b>	Valor líquido pago dos itens incluindo a taxa * Quantidade
<b>ss_net_profit</b>	<b>DECIMAL</b>	Lucro sobre o item comprador

#### 4.2.2 Volume de Dados

O volume de dados utilizado foi de 1 gigabyte tendo em vista as limitações da simulação de um *cluster* com apenas uma máquina. Além disso, com o objetivo de entender e apresentar o processo de criação de um DW para um cenário Big Data, este volume de dados é suficiente (MAJDALANY, ).

O data generator do TPCx-BB permite que o volume de dados seja escalonado, podendo chegar à casa do terabyte, portanto, com um *cluster* mais potente disponível seria possível aplicar o mesmo processo.

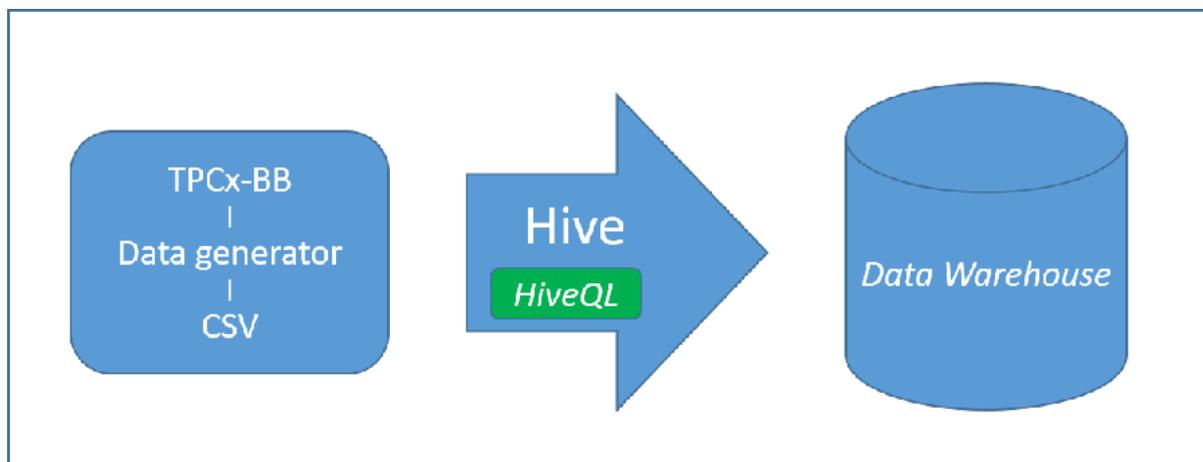
#### 4.2.3 Processo de ETL

Utilizando o *data generator* do TPCx-BB foi possível gerar os dados das respectivas tabelas num formato CSV<sup>1</sup>. O framework já realiza uma das etapas do ETL no que se refere a classificação e categorização dos dados semiestruturados e não-estruturados. Dessa forma, a arquitetura do DW poderia ser vista conforme ilustrada na Figura 12. Os dados gerados pelo framework da TPCx-BB foram utilizados como fonte de dados, e a linguagem HiveQL foi utilizada para carregar os dados no DW. Uma vez

<sup>1</sup> Comma-separated values; arquivos texto com células separadas por delimitadores.

que os dados foram gerados, foi necessário desenvolver os *scripts* para criação das tabelas usando a linguagem de consulta do Hive.

Figura 12 – Arquitetura do DW - TPCx-BB & Hive



Fonte: Elaborada pelo Autor (2016)

Com a *HiveQL* é possível realizar a carga das tabelas diretamente de arquivos CSV ao determinar qual o delimitador que será utilizado. Dessa forma, o processo de carga se dá de maneira simples.

#### 4.3 Comentários Finais

O processo de desenvolvimento de um DW para um cenário *Big Data* descrito neste capítulo utilizou ferramentas que auxiliam e facilitam bastante a evolução do projeto. Tecnologias como a CDH e o *Hive* se mostram bastante úteis ao tornar o desenvolvimento muito mais simples e próximo da realidade dos projetos de BI tradicionais, uma vez que a Cloudera permite transparência para a complexidade existente na instalação e compreensão do *Hadoop* e seu ecossistema, e o *Hive* acomoda usuários de SQL por meio das suas semelhanças.

## 5 Conclusão

Apesar de já se ter um grande acesso às informações e conseguir realizar processos de tomada de decisão por meio da utilização de técnicas de *Data Warehousing* sobre as tradicionais bases de dados, com o advento do *Big Data* torna-se clara a necessidade das empresas em expandir o conhecimento sobre o mercado e seus *stakeholders*.

Desse modo, a utilização de ferramentas como o *Hadoop* pode ser vista como uma tendência e pode ser um fator muito influente no processo de adaptação das empresas para este novo cenário.

O estudo sobre as tecnologias de *Big Data* é essencial para permitir determinar quais dessas tecnologias serão úteis para solucionar o problema de processamento, armazenamento e gerenciamento dos dados relacionados à empresa.

Neste trabalho, as tecnologias selecionadas trouxeram bastante aprendizado de modo a permitir a compreensão sobre o processamento e armazenamento de *Big Data*. Além disso, o desenvolvimento do DW, com o auxílio da tecnologia *Hive*, possibilitou a imersão nesse novo cenário de dados de modo a apresentar parte do potencial que esse tipo de aplicação possibilita, já que permite trabalhar com informações de diversas fontes, formatos e volumes, cada vez mais presentes na realidade empresarial.

### 5.1 Limitações

Devido ao desenvolvimento deste trabalho ter sido realizado por meio da simulação de um *cluster* utilizando uma só máquina (modo pseudodistribuído da CDH), algumas limitações de processamento fazem com que não se possa realizar testes completos sobre as tecnologias utilizadas.

Por esse motivo, o volume de dados utilizado – apesar de ser um volume suficiente para compreensão sobre o processo de desenvolvimento de um DW para *Big Data* – não foi escalonado para maiores quantidades. O *framework* de geração de dados da TPCx-BB possui um fator de escalonamento que pode ser definido e possibilita que o volume total de dados alcance a casa do *terabyte*, no entanto essa função não pôde ser devidamente explorada por limitações de *hardware*.

### 5.2 Trabalhos Futuros

Como foi observado neste trabalho, além do *Hive*, existem diversas tecnologias que compõem o ecossistema do *Hadoop*. O *Spark* é uma delas. Por meio dele também

seria possível realizar o desenvolvimento de um DW para *Big Data*, e assim seria possível analisar uma outra abordagem.

Além de avaliar outras tecnologias, também se pode ter como proposta de trabalho futuro a realização do mesmo processo sobre a base de dados gerada pela TPCx-BB, mas com um maior volume de dados. É possível, ainda, explorar o funcionamento do DW em *cluster* com um número significativo de máquinas a fim de analisar questões de desempenho do sistema.

Aplicar o processo desenvolvido neste trabalho num cenário real seria ideal para obter resultados mais concretos.

## Referências

- APACHE. *The Apache Software Foundation*. Disponível em: <<https://www.apache.org>>.
- BETTER, B. A. et al. Heading Towards Big Data. p. 220 – 225, 2013.
- CAMARGO, J. E. et al. A Big Data Analytics System to Analyze Citizens ' Perception of Security. p. 0 – 4, 2016.
- CUZZOCREA, A.; SACCÀ, D.; ULLMAN, J. D. Big data: A research agenda. *Proceedings of the 17th International Database Engineering {&} Applications Symposium on - IDEAS '13*, n. October, p. 198 – 203, 2013. Disponível em: <<http://dl.acm.org/citation.cfm?doid=2513591.2527071>>.
- EDUCAÇÃO, P. *METODOLOGIA CIENTÍFICA: TIPOS DE PESQUISA*. 2013. Disponível em: <<https://www.portaleducacao.com.br/pedagogia/artigos/50264/metodologia-cientifica-tipos-de-esquisa>>.
- INMON, W.; STRAUSS, D.; NEUSHLOSS, G. *Data Warehouse 2.0*. Morgan-Kaufmann. Disponível em: <<http://gen.lib.rus.ec/book/index.php?md5=4FACB22329B131FFEF49B85C37712225>>.
- JURNEY, R. *Agile Data Science: Building Data Analytics Applications with Hadoop*. [S.l.: s.n.], 2013. ISSN 1098-6596. ISBN 9781449326265.
- MAJDALANY, M. *Transaction Processing Performance Council*. Disponível em: <<http://www.tpc.org/>>.
- MELOROSE, J.; PERROY, R.; CAREAS, S. *Data Warehouse in the Age of Big Data*. [S.l.: s.n.], 2015. ISSN 1098-6596. ISBN 9788578110796.
- PROVOST, T. F. F. *Data Science for Business: What you need to know about data mining and data-analytic thinking*. [S.l.]: O'Reilly Media, 2013. ISBN 1449361323,9781449361327.
- SITTO, K.; PRESSER, M. *Field Guide to Hadoop*. [S.l.: s.n.], 2015. ISBN 0471462128.
- SUN, L. et al. Present situation and prospect of data warehouse architecture under the background of big data. *Proceedings - 2013 International Conference on Information Science and Cloud Computing Companion, ISCC-C 2013*, p. 529 – 535, 2014.
- THARAKAN, R. *What is scalability?* 2007. Disponível em: <<http://www.royans.net/wp/2007/09/22/what-is-scalability/>>.
- WHITE, T. *Hadoop : the definitive guide*. 2nd ed. ed. [S.l.]: O'Reilly, 2010. ISBN 9781449389734,1449389732.
- ZHANG, R. et al. Big Data for Medical Image Analysis: A Performance Study. *2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, p. 1660 – 1664, 2016. Disponível em: <<http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=7530064>>.