

SKRIPSI

REDUKSI BIG DATA DENGAN ALGORITMA CLUSTERING AGGLOMERATIVE UNTUK SISTEM TERDISTRIBUSI SPARK



Matthew Ariel

NPM: 2015730010

**PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS
UNIVERSITAS KATOLIK PARAHYANGAN
2019**

UNDERGRADUATE THESIS

**BIG DATA REDUCTION USING AGGLOMERATIVE CLUSTERING
ALGORITHM FOR SPARK DISTRIBUTED SYSTEM**



Matthew Ariel

NPM: 2015730010

**DEPARTMENT OF INFORMATICS
FACULTY OF INFORMATION TECHNOLOGY AND SCIENCES
PARAHYANGAN CATHOLIC UNIVERSITY
2019**

LEMBAR PENGESAHAN

REDUKSI BIG DATA DENGAN ALGORITMA CLUSTERING AGGLOMERATIVE UNTUK SISTEM TERDISTRIBUSI SPARK

Matthew Ariel

NPM: 2015730010

Bandung, 10 Desember 2019

Menyetujui,

Pembimbing

Dr. Veronica Sri Moertini

Ketua Tim Penguji

Anggota Tim Penguji

Mariskha Tri Adithia, P.D.Eng

Kristopher David Harjono, M.T.

Mengetahui,

Ketua Program Studi

Mariskha Tri Adithia, P.D.Eng

PERNYATAAN

Dengan ini saya yang bertandatangan di bawah ini menyatakan bahwa skripsi dengan judul:

REDUKSI BIG DATA DENGAN ALGORITMA CLUSTERING AGGLOMERATIVE UNTUK SISTEM TERDISTRIBUSI SPARK

adalah benar-benar karya saya sendiri, dan saya tidak melakukan penjiplakan atau pengutipan dengan cara-cara yang tidak sesuai dengan etika keilmuan yang berlaku dalam masyarakat keilmuan.

Atas pernyataan ini, saya siap menanggung segala risiko dan sanksi yang dijatuhkan kepada saya, apabila di kemudian hari ditemukan adanya pelanggaran terhadap etika keilmuan dalam karya saya, atau jika ada tuntutan formal atau non-formal dari pihak lain berkaitan dengan keaslian karya saya ini.

Dinyatakan di Bandung,
Tanggal 10 Desember 2019

Meterai
Rp. 6000

Matthew Ariel
NPM: 2015730010

ABSTRAK

Big data adalah istilah yang menggambarkan kumpulan data dalam jumlah yang sangat besar, baik data yang terstruktur maupun data yang tidak terstruktur. Kumpulan data tersebut menyimpan informasi yang bisa dianalisis dan diproses untuk memberikan wawasan kepada organisasi atau perusahaan. *Big data* dapat mencapai *petabyte* dan menghabiskan banyak tempat penyimpanan.

Big data perlu direduksi untuk menghemat tempat penyimpanan. Algoritma *Hierarchical Agglomerative Clustering* dapat digunakan untuk mereduksi data. Dengan bantuan sistem terdistribusi seperti Hadoop, proses reduksi data dapat dilakukan secara paralel dan lebih cepat. Sayangnya, teknologi Hadoop masih dapat dikatakan 'terlalu lambat' dalam melakukan proses reduksi data karena hasil sementara dari setiap tahap akan disimpan di *disk* sampai dibutuhkan kembali di tahap selanjutnya.

Untuk mempercepat proses reduksi data, Hadoop dapat digantikan dengan Spark. Spark adalah sistem terdistribusi, mirip seperti Hadoop. Tetapi, yang membedakan antara Hadoop dengan Spark adalah pada cara penyimpanan sementara saat melakukan proses reduksi data. Hadoop menggunakan *disk* sebagai tempat penyimpanan sementaranya, sedangkan Spark menggunakan memori sebagai tempat penyimpanan sementaranya. Pembacaan dan penulisan akan lebih cepat saat menggunakan memori dibandingkan dengan menggunakan *disk*, sehingga Spark akan lebih cepat dibandingkan dengan Hadoop.

Perangkat lunak dibuat untuk mengimplementasikan algoritma *Hierarchical Agglomerative Clustering* dalam Spark. Pengujian juga dilakukan dengan membandingkan waktu eksekusi algoritma *Hierarchical Agglomerative Clustering* saat diimplementasikan pada Hadoop dan saat diimplementasikan pada Spark. Waktu eksekusi dicatat untuk ukuran data 1GB, 2GB, 3GB, 5GB, 10GB, 15GB, dan 20GB.

Berdasarkan hasil pengujian, Spark memiliki waktu eksekusi yang lebih cepat dibandingkan dengan Hadoop pada jumlah partisi yang besar. Waktu eksekusi Spark menurun ketika jumlah partisi ditingkatkan, sedangkan waktu eksekusi Hadoop menurun ketika jumlah partisi ditingkatkan. Waktu eksekusi terbaik Spark masih lebih cepat dibandingkan waktu eksekusi terbaik Hadoop.

Kata-kata kunci: *Big Data*, Reduksi Data, *Hierarchical Agglomerative Clustering*, Spark, Hadoop

ABSTRACT

Big data is a term that describes the large volume of data, both structured and unstructured. The data set stores information can be analyzed and processed to provide insight to organization or company. Big data can reach up to petabytes and takes a lot of storage spaces.

Big data need to be reduce to save storage space. The Hierarchical Agglomerative Clustering algorithm can be used to reduce data. With the help of distributed systems such as Hadoop, reduction process can be done in parallel with less execution time. Unfortunately, Hadoop can still be said to be 'too slow' in the process of data reduction because temporary results from each stage will be stored on the disk until it is needed again at a later stage.

To speed up the data reduction process, Hadoop can be replaced with Spark. Spark is a distributed system, similar to Hadoop. However, what distinguishes Hadoop from Spark is the way Spark temporarily store data. Hadoop uses disk as its temporary storage, while Spark uses memory as its temporary storage. Read and write process will be faster when using memory than using disks, Spark will be faster than Hadoop.

The Hierarchical Agglomerative Clustering algorithm is implemented in the software. Experiment were done by comparing the execution time of the Hierarchical Agglomerative Clustering algorithm when implemented on Hadoop and Spark. The execution time is recorded for 1GB, 2GB, 3GB, 5GB, 10GB, 15GB, dan 20GB of data.

Based on the experiment, Spark has a faster execution time compared to Hadoop on a large number of partitions. Spark execution time decreases when the number of partitions is increased, whereas Hadoop execution time decreases when the number of partitions is increased. Spark best executiron time is still much better than Hadoop best execution time.

Keywords: Big Data, Data Reduction, Hierarchical Agglomerative Clustering, Spark, Hadoop

Dipersembahkan untuk Tuhan YME, keluarga tercinta, Ibu Veronica sebagai dosen pembimbing, teman-teman yang telah berperan dalam pembuatan skripsi ini, dan diri sendiri.

KATA PENGANTAR

Puji dan syukur kehadirat Tuhan Yang Maha Esa atas berkat rahmat serta kasih-Nya sehingga penulis dapat menyelesaikan skripsi ini yang mengambil judul Reduksi Big Data dengan Algoritma Clustering Agglomerative untuk Sistem Terdistribusi Spark. Penulisan skripsi ini diajukan untuk memenuhi salah satu syarat untuk memperoleh gelar Sarjana pada Program Studi Teknik Informatika Universitas Katolik Parahyangan. Pada penyusunan dan penulisan skripsi ini, penulis menyadari bahwa penyusunan dan penulisan skripsi ini juga tidak terlepas dari bantuan berbagai pihak, baik langsung maupun tidak langsung. Secara khusus, penulis ingin berterima kasih kepada:

1. Keluarga yang selalu memberi dukungan dan mengurus penulis selama mengerjakan skripsi.
2. Ibu Veronica Sri Moertini selaku dosen pembimbing yang telah membimbing penulis dan memberikan dukungan dan bantuan kepada penulis dalam proses penyusunan skripsi ini.
3. Sahabat yang selalu memberi dukungan dan semangat, khusunya Edick Zakari, Felicia Christiany, Irvan Wijaya, Emanuel Yudistira, Hapsari Laksmi, Kezia, Raymond Nagawijaya, Cornelius David, Aditya Putra, Thoby, dan teman-teman yang tidak dapat saya sebutkan.
4. Sahabat seperjuangan dari jurusan lain Rensky Picco dan Christian Stefano.
5. Sahabat SMURF -rm -rf, Wanted To Resign, sad, Mati Lari Sudah, dan shitos.
6. Teman dan kerabat lainnya yang tidak bisa saya sebutkan satu-persatu.

Penulis menyadari bahwa skripsi ini masih jauh dari kata sempurna. Oleh karena itu, penulis memohon maaf jika terdapat kesalahan. Penulis juga mengharapkan kritik dan saran yang membangun untuk menyempurnakan skripsi ini. Semoga skripsi ini dapat memberi informasi yang bermanfaat dan menjadi inspirasi untuk penelitian-penelitian berikutnya.

Bandung, Desember 2019

Penulis

DAFTAR ISI

KATA PENGANTAR	xv
DAFTAR ISI	xvii
DAFTAR GAMBAR	xix
DAFTAR TABEL	xxiii
1 PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	2
1.3 Tujuan	2
1.4 Batasan Masalah	2
1.5 Metodologi	2
1.6 Sistematika Pembahasan	3
2 LANDASAN TEORI	5
2.1 <i>Big Data</i>	5
2.2 Algoritma Hierarchical Clustering	6
2.3 Hadoop	12
2.3.1 Hadoop Distributed File System (HDFS)	13
2.3.2 MapReduce [5]	15
2.3.3 YARN	16
2.3.4 Pembahasan Algoritma <i>Agglomerative Clustering</i> pada Hadoop [1]	17
2.4 Spark [7]	18
2.4.1 Komponen Spark [7]	18
2.4.2 Tiga Cara Membangun Spark di Atas Hadoop	20
2.4.3 Arsitektur Spark	20
2.4.4 <i>Resilient Distributed Datasets (RDD)</i>	21
2.5 Scala	24
2.5.1 <i>Expressions</i>	24
2.5.2 <i>Blocks</i>	25
2.5.3 <i>Loop</i> dan <i>Conditional</i>	25
2.5.4 <i>Functions</i>	26
2.5.5 <i>Methods</i>	26
2.5.6 <i>Class</i> dan <i>Object</i>	27
2.5.7 <i>Higher Order Function</i>	28
3 STUDI DAN EKSPLORASI APACHE SPARK	31
3.1 Instalasi Apache Spark	31
3.2 Eksplorasi Spark Shell	32
3.3 Instalasi Apache Spark pada <i>Multi-Node Cluster</i>	34
3.4 Percobaan Spark Submit	36

4 ANALISIS DAN PERANCANGAN	43
4.1 Analisis Masalah	43
4.1.1 Identifikasi Masalah	43
4.1.2 Analisis Masukan dan Keluaran	44
4.1.3 Analisis <i>Hierarchical Agglomerative Clustering MapReduce</i>	45
4.1.4 Diagram Alur	47
4.1.5 Analisis <i>Hierarchical Agglomerative Clustering</i> pada Spark	52
4.2 Perancangan Perangkat Lunak	58
4.2.1 Diagram <i>Use Case</i> dan Skenario	58
4.2.2 Diagram Kelas	60
4.2.3 Rancangan Antarmuka	63
5 IMPLEMENTASI DAN PENGUJIAN PERANGKAT LUNAK	69
5.1 Implementasi Perangkat Lunak	69
5.1.1 Lingkungan Perangkat Kerat	69
5.1.2 Lingkungan Perangkat Lunak	69
5.1.3 <i>User Interface</i>	70
5.2 Pengujian Fungsional Perangkat Lunak	72
5.3 Pengujian Eksperimental Perangkat Lunak	74
5.3.1 Pengujian Dampak Partisi Pada Waktu Eksekusi	75
5.3.2 Pengujian Dengan Batas Objek Maksimum yang Berbeda	85
5.3.3 Pengujian Dengan Metode Linkage yang Berbeda	93
6 KESIMPULAN DAN SARAN	97
6.1 Kesimpulan	97
6.2 Saran	97
DAFTAR REFERENSI	99
A KODE PROGRAM SPARK	101
B KODE PROGRAM UNTUK ANTARMUKA	107
C KODE PROGRAM MAPREDUCE	113

DAFTAR GAMBAR

2.1	Karakteristik <i>big data</i>	5
2.2	Matriks jarak	6
2.3	Matriks jarak	7
2.4	<i>dendrogram</i>	7
2.5	Metode <i>single linkage</i>	8
2.6	Metode <i>complete linkage</i>	8
2.7	Metode <i>centroid linkage</i>	9
2.8	Matriks jarak	9
2.9	Hasil penggabungan <i>cluster</i>	10
2.10	Hasil rekalkulasi	11
2.11	Hasil akhir <i>dendrogram</i>	11
2.12	Perpotongan <i>dendrogram</i>	11
2.13	Modul-modul Hadoop	13
2.14	Arsitektur HDFS	14
2.15	Arsitektur MapReduce	15
2.16	Proses MapReduce	16
2.17	Proses menjalankan aplikasi pada YARN	17
2.18	Komponen pada Spark	18
2.19	Macam-macam cara instalasi Spark	20
2.20	Arsitektur Spark	20
3.1	<i>Spark Shell</i>	32
3.2	<i>Word Count</i>	33
3.3	IntelliJ IDEA	36
3.4	Proyek sbt	37
3.5	Konfigurasi proyek	38
3.6	Struktur proyek	38
3.7	Konfigurasi sbt	39
3.8	<i>object WordCount</i>	39
3.9	JAR	40
3.10	Hasil perintah 'sbt package'	40
3.11	Penggumpulan JAR kepada <i>spark-submit</i>	40
3.12	Alamat Spark UI	41
3.13	Spark UI	41
4.1	Penulisan kepada disk di MapReduce	44
4.2	Penulisan kepada memori di Spark	44
4.3	Contoh <i>Input</i> Program	44
4.4	Contoh <i>Output</i> Program	45
4.5	Diagram alur perangkat lunak	47
4.6	Partisi RDD	48
4.7	Contoh Perubahan <i>Block</i> HDFS menjadi RDD <i>String</i>	48
4.8	RDD <i>parsing</i> dan kelas <i>Node</i>	49

4.9	<i>Worker</i> memproses partisi	49
4.10	Contoh pemasangan nilai acak kepada objek <i>Node</i>	49
4.11	Pengelompokkan <i>Node</i> berdasarkan <i>key</i>	50
4.12	Proses reduksi dan kelas <i>Pattern</i>	50
4.13	Contoh Proses Komputasi Pola	51
4.14	Perpotongan <i>dendrogram</i>	51
4.15	Penyimpanan pola pada HDFS	52
4.16	Contoh perhitungan matriks dan pembentukan dendrogram	57
4.17	Contoh pemotongan <i>dendrogram</i>	57
4.18	Diagram <i>use case</i> perangkat lunak <i>Hierarchical Agglomerative Clustering</i>	59
4.19	Diagram kelas	60
4.20	Kelas Main, SparkConfig, SparkContext	60
4.21	Kelas DataReducer	61
4.22	Kelas Dendrogram	61
4.23	Kelas Cluster	62
4.24	Kelas Pattern	62
4.25	Kelas Node	63
4.26	Rancangan antarmuka menu Jalankan Program	64
4.27	Halaman web Hadoop	65
4.28	Rancangan antarmuka menu Lihat Pola	65
4.29	Rancangan antarmuka halaman partisi	66
4.30	Rancangan antarmuka halaman pola	66
4.31	Halaman web HDFS	67
5.1	Tampilan menu <i>Submit</i>	70
5.2	Tampilan menu <i>Data</i>	70
5.3	Tampilan halaman <i>list</i>	71
5.4	Tampilan halaman data	72
5.5	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 1GB, jumlah objek maksimum 30, dan total 10 core	76
5.6	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 2GB, jumlah objek maksimum 30, dan total 10 core	77
5.7	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 3GB, Objek Maksimum 30, dan Total 10 Core	78
5.8	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 5GB, Objek Maksimum 30, dan Total 10 Core	79
5.9	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 10GB, Objek Maksimum 30, dan Total 10 Core	81
5.10	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 15GB, Objek Maksimum 30, dan Total 10 Core	83
5.11	Hasil Percobaan Partisi Spark dan Hadoop dengan ukuran data 20GB, Objek Maksimum 30, dan Total 10 Core	84
5.12	Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan ukuran data 5GB, Objek Maksimum 50, dan Total 10 Core	86
5.13	Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 10 GB, Objek Maksimum 50, dan Total 10 Core	87
5.14	Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 15 GB, Objek Maksimum 50, dan Total 10 Core	88
5.15	Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 20 GB, Objek Maksimum 50, dan Total 10 Core	89
5.16	Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 5 GB, Objek Maksimum 100, dan Total 10 Core	90

5.17 Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 10 GB, Objek Maksimum 100, dan Total 10 Core	91
5.18 Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 15 GB, Objek Maksimum 100, dan Total 10 Core	92
5.19 Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 20 GB, Objek Maksimum 100, dan Total 10 Core	93
5.20 Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 10 GB, Objek Maksimum 30, dan Total 10 Core	94
5.21 Hasil Percobaan Jumlah Partisi Spark dan Hadoop dengan Ukuran Data 10 GB, Objek Maksimum 30, dan Total 10 Core	95

DAFTAR TABEL

2.1	Tabel Data Koordinat	9
2.2	Tabel Contoh Data <i>Cluster</i>	12
2.3	Tabel Hasil Pola Cluster A	12
2.4	Tabel transformations	23
2.5	Tabel Actions	24
4.1	Tabel Contoh Data <i>Cluster</i>	52
4.2	Tabel Hasil Pola Cluster A	52
5.1	Tabel data yang digunakan pada eksperimen	75
5.2	Percobaan Jumlah Partisi Hadoop dan Spark dengan Ukuran Data 1 GB	76
5.3	Percobaan Jumlah Partisi Hadoop dan Spark dengan Ukuran Data 2 GB	77
5.4	Percobaan Jumlah Partisi Hadoop dan Spark dengan Ukuran Data 3 GB	78
5.5	Percobaan Jumlah Partisi Hadoop dan Spark dengan Ukuran Data 5 GB	79
5.6	Percobaan Jumlah Partisi Spark dengan Ukuran Data 10 GB	80
5.7	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 10 GB	80
5.8	Percobaan Jumlah Partisi Spark dengan Ukuran Data 15 GB	82
5.9	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 15 GB	82
5.10	Percobaan Jumlah Partisi Spark dengan Ukuran Data 20 GB	84
5.11	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 20 GB	84
5.12	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 5 GB	85
5.13	Percobaan Jumlah Partisi Spark dengan Ukuran Data 5 GB	85
5.14	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 10 GB	86
5.15	Percobaan Jumlah Partisi Spark dengan Ukuran Data 10 GB	86
5.16	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 15 GB	87
5.17	Percobaan Jumlah Partisi Spark dengan Ukuran Data 15 GB	87
5.18	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 20 GB	88
5.19	Percobaan Jumlah Partisi Spark dengan Ukuran Data 20 GB	88
5.20	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 5 GB	89
5.21	Percobaan Jumlah Partisi Spark dengan Ukuran Data 5 GB	89
5.22	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 10 GB	90
5.23	Percobaan Jumlah Partisi Spark dengan Ukuran Data 10 GB	90
5.24	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 15 GB	91
5.25	Percobaan Jumlah Partisi Spark dengan Ukuran Data 15 GB	91
5.26	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 20 GB	92
5.27	Percobaan Jumlah Partisi Spark dengan Ukuran Data 20 GB	92
5.28	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 10 GB	93
5.29	Percobaan Jumlah Partisi Spark dengan Ukuran Data 10 GB	94
5.30	Percobaan Jumlah Partisi Hadoop dengan Ukuran Data 10 GB	94
5.31	Percobaan Jumlah Partisi Spark dengan Ukuran Data 10 GB	95

BAB 1

PENDAHULUAN

1.1 Latar Belakang

Big data adalah sebuah istilah yang menggambarkan volume data yang besar, baik data yang terstruktur maupun data yang tidak terstruktur. Data-data tersebut memiliki potensi untuk digali menjadi informasi yang penting. Dalam bidang *big data* ada beberapa tantangan seperti volume data yang besar, kecepatan aliran data yang masuk, dan variasi data dengan format yang berbeda. Tantangan tersebut membuat aplikasi pemrosesan data tradisional tidak bisa memproses dan menganalisis *big data*. Muncul teknologi-teknologi seperti Hadoop dan Spark yang dirancang khusus untuk menangani *big data*.

Big data akan lebih mudah dianalisis dan diterapkan teknik-teknik *data-mining* ketika volume *big data* tersebut telah direduksi. Reduksi data adalah proses mengecilkan ukuran data dengan mengambil rangkuman dari sekelompok data. *Big data* akan dipecah menjadi beberapa bagian. Dari bagian tersebut akan diolah menggunakan algoritma *Hierarchical Agglomerative Clustering* untuk menghasilkan *cluster-cluster*. Dari *cluster-cluster* yang dihasilkan, akan diambil rangkuman yang bisa menggambarkan karakteristik setiap *cluster*. Rangkuman yang menggambarkan karakteristik *cluster* disebut sebagai pola. Pola mengandung jumlah total objek, nilai minimum, maksimum, nilai rata-rata, dan standar-deviasi dari setiap atribut pada *cluster*. Nilai-nilai tersebut diambil karena dapat menggambarkan karakteristik *cluster*. Dengan mereduksi data, kita bisa menghemat biaya pengiriman data, *disk space*, dan jumlah data yang diproses. Hasil dari reduksi *big data* harus bisa mewakili data yang belum direduksi secara akurat.

Algoritma *Hierarchical Agglomerative Clustering* membangun hierarki dari sekelompok data. Setiap objek akan ditempatkan kepada *cluster*-nya tersendiri pada awalnya. Kemudian, *cluster* terdekat akan digabung menjadi satu *cluster*. Penggabungan *cluster* terdekat akan diulang sampai hanya satu *cluster* yang tersisa. Proses pembangunan hierarki dapat digambarkan dengan *dendrogram*. *Dendrogram* akan dipotong untuk menghasilkan *cluster-cluster*. Dari setiap *cluster* hasil perpotongan, akan dicari polanya dan disimpan sebagai hasil akhir.

Untuk mempercepat proses reduksi, pekerjaan dapat dipecah dan dikerjakan secara paralel dengan bantuan Hadoop. Dengan bantuan Hadoop, proses reduksi data akan lebih cepat. Implementasi algoritma *Hierarchical Agglomerative Clustering* sudah dilakukan pada Hadoop [1]. Hasil dari penelitian tersebut membuktikan bahwa algoritma *Hierarchical Agglomerative Clustering* sudah dapat mereduksi data dengan menyimpan pola hasil reduksi. Tetapi ada beberapa limitasi yang dimiliki Hadoop.

Walau Hadoop dapat melakukan proses reduksi secara paralel, waktu yang dibutuhkan Hadoop untuk melakukan reduksi data masih terlalu lambat. Hadoop banyak melakukan proses penulisan dan pembacaan kepada disk. Dari satu tahap ke tahap lainnya, Hadoop akan menulis dan membaca hasil sementara kepada disk. Hadoop perlu digantikan dengan Spark untuk mencepat proses reduksi.

Spark adalah *distributed cluster-computing framework* yang bisa menggantikan MapReduce beserta kekurangannya. *In-memory processing* pada Spark dapat mengalahkan kecepatan pemrosesan pada Hadoop MapReduce. Karena data disimpan pada RAM, kecepatan pemrosesan akan jauh lebih cepat. Spark membaca data yang akan direduksi dari RAM. Pembacaan data dari RAM akan lebih cepat dibanding disk.

Pada skripsi ini, dibangun sebuah perangkat lunak yang dapat mereduksi *big data*. Perangkat lunak tersebut akan dibangun menggunakan *framework* terdistribusi Spark dan mengimplementasikan algoritma *Hierarchical Agglomerative Clustering* yang khusus dirancang untuk lingkungan Spark. Perangkat lunak

dapat menampilkan hasil reduksi dalam format tabel. Dengan menggunakan Spark, waktu proses reduksi data menjadi lebih cepat dibanding MapReduce.

1.2 Rumusan Masalah

Berdasarkan latar belakang di atas, dapat dibentuk rumusan masalah sebagai berikut:

1. Bagaimana cara kerja algoritma *Hierarchical Agglomerative Clustering* berbasis MapReduce untuk mereduksi *big data*?
2. Bagaimana cara mengkustomisasi dan mengimplementasikan algoritma *Agglomerative Clustering* pada sistem tersebut Spark?
3. Bagaimana mengukur performa hasil dari implementasi algoritma *Agglomerative Clustering* pada sistem tersebut Spark?

1.3 Tujuan

Berdasarkan rumusan masalah di atas, tujuan dari penelitian adalah sebagai berikut:

1. Mempelajari cara kerja algoritma *Hierarchical Agglomerative Clustering* berbasis MapReduce untuk mereduksi *big data*.
2. Mengkustomisasi dan mengimplementasikan algoritma *Hierarchical Agglomerative Clustering* pada lingkungan Spark.
3. Melakukan eksperimen pada Spark dan Hadoop untuk membandingkan waktu eksekusi kedua perangkat lunak.

1.4 Batasan Masalah

Batasan masalah pada skripsi ini adalah sebagai berikut:

1. Pola yang dikomputasi hanya sebatas jumlah objek, nilai minimum, maksimum, rata-rata dan standar deviasi.
2. Perangkat lunak yang dibangun tidak menangani pembersihan masukkan *big data*.

1.5 Metodologi

Metodologi yang digunakan dalam pembuatan skripsi ini adalah:

1. Melakukan studi literatur Hadoop hanya mempelajari konsep dasar dari Hadoop dan *Hadoop Distributed File System* (HDFS).
2. Melakukan studi literatur tentang konsep Apache Spark.
3. Melakukan studi literatur bahasa pemrograman Scala.
4. Melakukan studi literatur tentang algoritma *Hierarchical Agglomerative Clustering*.
5. Melakukan instalasi dan konfigurasi Apache Spark.
6. Melakukan eksperimen dengan bahasa pemrograman Scala.
7. Melakukan eksperimen dengan Spark RDD.

8. Melakukan kustomisasi algoritma *Hierarchical Agglomerative Clustering* untuk Spark.
9. Mencari dan mengumpulkan data uji coba yang bervolume besar.
10. Merancang dan mengimplementasikan perangkat lunak.
11. Melakukan eksperimen terhadap perangkat lunak dan menganalisis hasil eksperimen.
12. Menulis dokumen skripsi.

1.6 Sistematika Pembahasan

Laporan penelitian tersusun ke dalam enam bab secara sistematis sebagai berikut:

- Bab 1 Pendahuluan
Berisi latar belakang, rumusan masalah, tujuan, batasan masalah, metodologi penelitian, dan sistematika pembahasan.
- Bab 2 Dasar Teori
Berisi dasar teori tentang *big data*, *Hierarchical Agglomerative Clustering*, Hadoop, Spark, dan Scala.
- Bab 3 Studi dan Eksplorasi Apache Spark
Berisi percobaan-percobaan yang dilakukan pada Spark.
- Bab 4 Analisis dan Perancangan
Berisi analisis masalah, diagram alur, *use case* dan skenario, diagram kelas, dan perancangan antarmuka.
- Bab 5 Implementasi dan Pengujian
Berisi implementasi antarmuka perangkat lunak, pengujian eksperimen, dan kesimpulan dari pengujian.
- Bab 5 Implementasi dan Pengujian
Berisi kesimpulan awal sampai akhir penelitian dan saran untuk penelitian selanjutnya.