

**SKRIPSI**

**ANALISIS CO-OCCURRENCE NETWORKS PADA BIG GRAPH**



**Indra Permana Sugianto**

**NPM: 2017730008**

**PROGRAM STUDI TEKNIK INFORMATIKA  
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS  
UNIVERSITAS KATOLIK PARAHYANGAN  
2022**

# LEMBAR PENGESAHAN

## ANALISIS CO-OCCURRENCE NETWORKS PADA BIG GRAPH

Indra Permana Sugianto

NPM: 2017730008

Bandung, 20 Januari 2022

Menyetujui,

Pembimbing

Digitally signed  
by Veronica Sri  
Moertini

Dr. Veronica Sri Moertini

Ketua Tim Penguji

Digitally signed  
by Luciana  
Abednego

Luciana Abednego, M.T.

Anggota Tim Penguji

Digitally signed  
by Lionov

Lionov, Ph.D.

Mengetahui,

Ketua Program Studi

Digitally signed  
by Mariskha Tri  
Adithia

Mariskha Tri Adithia, P.D.Eng

## PERNYATAAN

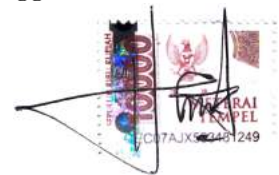
Dengan ini saya yang bertandatangan di bawah ini menyatakan bahwa skripsi dengan judul:

### **ANALISIS CO-OCCURRENCE NETWORKS PADA BIG GRAPH**

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 20 Januari 2022

A handwritten signature in black ink is written over a 1000 Rupiah postage stamp. The stamp features the Garuda Pancasila emblem and the text 'REPUBLIK INDONESIA', '1000', and 'POSTAL SERVICE'. The signature is written in a cursive style.

Indra Permana Sugianto  
NPM: 2017730008

## ABSTRAK

Salah satu tren analisis data saat ini adalah analisis graf dari data media sosial. Data media sosial merupakan contoh *big data* yang dapat disiapkan menjadi graf dan dianalisis menggunakan teknologi *big data*. Permasalahan seperti deteksi komunitas dan analisis profil komunitas, menjadi salah satu topik yang banyak diteliti saat ini. Berbagai algoritma deteksi komunitas, seperti *Connected Component*, *Strongly Connected Component*, dan *Triangle Count* telah diimplementasikan. Namun, sebagian besar masih untuk graf tidak berarah, atau membutuhkan input jumlah komunitas, atau banyak algoritma yang komputasinya lambat. Oleh karena itu, kurang cocok diimplementasikan untuk *big graph*, yaitu *big data* yang mendeskripsikan relasi antar objek.

Pada penelitian sebelumnya, berhasil diciptakan solusi untuk deteksi komunitas pada *directed big graph* dari data Twitter. Komunitas dapat didefinisikan sebagai sekumpulan pengguna yang aktif berinteraksi pada periode tertentu karena membahas topik tertentu. Solusi yang diusulkan dengan menggunakan *Motif Finding*, yaitu teknik pengenalan pola (*pattern matching*) yang diimplementasikan pada *library* GraphFrames di Apache Spark. Dilakukan perbandingan deteksi komunitas menggunakan algoritma *Strongly Connected Component* (SCC) dan teknik *Motif Finding*. Secara waktu eksekusi, didapatkan bahwa teknik *Motif Finding* memperoleh hasil yang lebih cepat dibandingkan SCC. Namun, teknik *Motif Finding*, perlu mendefinisikan pola-pola komunitas (motif) yang ingin dideteksi.

Skripsi ini berfokus pada deteksi komunitas dan analisis *Co-Occurrence* untuk analisis profil komunitas dari data Twitter. *Co-Occurrence* bermakna kemunculan suatu kejadian dalam frekuensi yang tinggi. Contohnya, kemunculan kata “*trading*” dalam kumpulan teks *tweet* di domain investasi. Deteksi komunitas menggunakan teknik *Motif Finding*, sedangkan untuk analisis profilnya menggunakan analisis *tweet* yang sering dikirimkan di komunitas (*co-occur*). Analisis *Co-Occurrence* dilakukan dengan teknik pemodelan topik.

Pengujian metode deteksi komunitas dan analisis profil komunitas dilakukan pada data *tweet* dengan domain Covid berbahasa Indonesia. Berdasarkan eksperimen yang dilakukan, metode di atas dapat mendeteksi komunitas dan menganalisis profilnya melalui kata-kata yang sering dikirimkan. Kata tersebut dapat digunakan untuk mendeskripsikan konten yang dibicarakan pada komunitas yang terbentuk. Namun, kekurangan dari metode di atas adalah harus mengetahui jumlah topiknya secara pasti, jika ingin hasilnya lebih baik.

**Kata-kata kunci:** *big graph*, deteksi komunitas, *co-occurrence*, *motif finding*, Apache Spark, GraphFrames, pemodelan topik, Twitter

## ABSTRACT

One of the current trends in data analysis is graph analysis of social media data. Social media data is an example of big data that can be prepared into graphs and analyzed using big data technologies. Problems such as community detection and analysis of community profiles (profiling) have become one of the most studied topics today. Various community detection algorithms, such as Connected Component, Strongly Connected Component, and Triangle Count have been implemented. However, most of them are still for undirected graphs, or require input number of communities, or many algorithms that are computationally slow. Therefore, it is not suitable to be implemented for big graphs, namely big data that describes the relation between objects.

In the previous research, a solution was created for community detection on the directed big graph of Twitter data. The community can be defined as a group of users who actively interact during a certain period, because they discuss certain topics. The proposed solution uses Motif Finding, a pattern matching technique that implemented in the GraphFrames library in Apache Spark. A comparison of community detection was carried out using the Strongly Connected Component (SCC) algorithm and the Motif Finding technique. In terms of execution time, it was found that the Motif Finding technique obtained faster results than SCC. However, the Motif Finding technique needs to define community patterns (motifs) to be detected.

This thesis focuses on community detection and Co-Occurrence analysis for community profile analysis from Twitter data. Co-Occurrence can be defined as the occurrence of an event in a high frequency. For example, the occurrence of the word “trading” in a collection of tweet text in the investment domain. The community detection uses the Motif Finding technique, while the profiling uses the analysis of tweets that are often sent in the community (co-occur). Co-Occurrence analysis was carried out using topic modeling techniques.

The evaluation of the community detection and profiling methods were carried out on tweet data with the Indonesian language Covid domain. Based on the experiments, the methods are able to detect communities and analyze the community profiles through frequently sent words. The words can be used to describe the content discussed in the community that is formed. However, the drawback of the methods are that necessary to know the number of topics, if expecting to get better results.

**Keywords:** big graph, community detection, co-occurrence, motif finding, Apache Spark, GraphFrames, topic modelling, Twitter

## KATA PENGANTAR

Puji syukur Penulis panjatkan ke hadirat Tuhan yang Maha Esa, atas berkat-Nya Penulis mampu menyelesaikan skripsi ini. Pada kesempatan ini, Penulis ingin berterima kasih kepada pihak yang membantu dan mendukung selama menjalani kuliah dan pengerjaan skripsi, yaitu sebagai berikut.

1. Kepada orang tua dan paman Penulis yang selalu memberi doa, dukungan dan semangat;
2. Kepada Ibu Dr. Veronica Sri Moertini selaku pembimbing. Terima kasih atas ilmu, pengalaman, semangat, dan bimbingan yang telah diberikan;
3. Kepada Ibu Luciana Abednego, M.T. dan Bapak Lionov, Ph.D. selaku penguji. Terima kasih atas masukan dan saran yang diberikan, sehingga skripsi ini dapat dikembangkan menjadi lebih baik;
4. Kepada seluruh dosen yang mengajar Penulis sejak pertama masuk kuliah. Terima kasih atas bekal ilmu dan pengalaman yang berharga;
5. Kepada teman-teman mahasiswa bimbingan Bu Vero (Fariz, Yoga, Fritz, Yohan), terima kasih sudah berjuang bersama. Akhirnya semua bisa lulus bareng :D;
6. Kepada Fanny Susiani Bunawan, terima kasih atas dukungan yang telah diberikan;
7. Kepada teman-teman mahasiswa angkatan 2017 sampai 2020. Terima kasih atas momen kebersamaannya. Semangat menjalani kuliahnya, bagi teman-teman adik tingkat;
8. Kepada seluruh pihak yang tidak Penulis cantumkan pada bagian ini. Terima kasih atas dukungannya.

Sebagai penutup, Penulis berterima kasih kepada pembaca dokumen skripsi ini. Penulis berharap semoga penelitian yang dilakukan dapat bermanfaat bagi para pembaca, kemajuan ilmu pengetahuan, dan dapat menjadi referensi yang baik bagi penelitian di masa depan. Penulis memohon maaf jika terdapat kesalahan penulisan, sistematika, atau metode yang dipaparkan pada skripsi ini. Penulis sangat terbuka untuk menerima masukan dan saran tambahan dari pembaca.

Bandung, Januari 2022

Penulis

# DAFTAR ISI

<b>KATA PENGANTAR</b>	<b>xv</b>
<b>DAFTAR ISI</b>	<b>xvii</b>
<b>DAFTAR GAMBAR</b>	<b>xix</b>
<b>DAFTAR KODE PROGRAM</b>	<b>xxiii</b>
<b>1 PENDAHULUAN</b>	<b>1</b>
1.1 Latar Belakang . . . . .	1
1.2 Rumusan Masalah . . . . .	3
1.3 Tujuan . . . . .	3
1.4 Batasan Penelitian . . . . .	3
1.5 Metodologi . . . . .	3
1.6 Sistematika Pembahasan . . . . .	4
<b>2 LANDASAN TEORI</b>	<b>5</b>
2.1 Big Data . . . . .	5
2.2 Graf [1] . . . . .	7
2.3 Co-Occurrence Networks . . . . .	8
2.4 Apache Spark . . . . .	8
2.4.1 Definisi dan Arsitektur Spark . . . . .	8
2.4.2 Konsep RDD dan Spark Dataframe . . . . .	9
2.4.3 Proses Komputasi Aplikasi Spark . . . . .	10
2.5 GraphX dan GraphFrames . . . . .	12
2.5.1 Pengantar . . . . .	12
2.5.2 Motif Finding . . . . .	13
2.6 Temporal Active Communities (TAC) . . . . .	13
2.7 Pemodelan Topik . . . . .	16
2.7.1 TF-IDF sebagai Bobot Fitur . . . . .	17
2.7.2 Latent Semantic Analysis untuk Pemodelan Topik . . . . .	17
<b>3 STUDI EKSPLORASI APACHE SPARK, GRAPHFRAMES, DAN SCIKIT-LEARN</b>	<b>21</b>
3.1 Eksperimen Analisis Co-Occurrence Networks Data Sitasi Indeks MEDLINE menggunakan GraphX [2] . . . . .	21
3.1.1 Pengantar . . . . .	21
3.1.2 Deskripsi Dataset . . . . .	21
3.1.3 Eksplorasi dan Penyiapan Data . . . . .	21
3.1.4 Pemodelan Graf . . . . .	23
3.1.5 Hasil Analisis . . . . .	23
3.2 Eksperimen Analisis Data Rute Sepeda di San Francisco menggunakan SparkSQL dan GraphFrames . . . . .	26
3.2.1 Pengantar . . . . .	26

3.2.2	Deskripsi Dataset . . . . .	26
3.2.3	Eksplorasi dan Penyiapan Data . . . . .	26
3.2.4	Pemodelan Graf dan Hasil Analisis . . . . .	29
3.3	Eksperimen Latent Semantic Analysis pada Library Scikit-Learn . . . . .	31
3.3.1	Pengantar . . . . .	31
3.3.2	Hasil Eksperimen . . . . .	31
3.4	Eksperimen Deteksi Komunitas dan Analisis Profil Komunitas dari Data Interaksi Quote Pengguna Twitter pada Domain Covid . . . . .	36
3.4.1	Pengantar . . . . .	36
3.4.2	Deskripsi Dataset . . . . .	36
3.4.3	Eksplorasi dan Penyiapan Data . . . . .	36
3.4.4	Pemodelan Graf . . . . .	38
3.4.5	Deteksi Komunitas . . . . .	39
3.4.6	Analisis Profil Komunitas . . . . .	40
<b>4</b>	<b>PENGUMPULAN, EKSPLORASI, DAN PENYIAPAN DATA</b>	<b>47</b>
4.1	Pengumpulan Data . . . . .	47
4.1.1	Arsitektur Sistem Pengumpul Data . . . . .	47
4.1.2	Deskripsi Format Data Twitter . . . . .	47
4.1.3	Deskripsi Data Studi Kasus . . . . .	49
4.2	Penyiapan Data . . . . .	49
4.3	Eksplorasi Data . . . . .	52
4.4	Pemodelan Graf . . . . .	56
<b>5</b>	<b>ANALISIS DATA DAN EVALUASI</b>	<b>57</b>
5.1	Analisis dan Interpretasi Hasil Analisis . . . . .	57
5.1.1	Deteksi Komunitas dan Analisis Profil Komunitas pada Data Domain Hidroponik . . . . .	57
5.1.2	Deteksi Komunitas dan Analisis Profil Komunitas pada Data Domain Kuliner Indonesia . . . . .	66
5.1.3	Deteksi Komunitas dan Analisis Profil Komunitas pada Data Domain Pemasaran Global . . . . .	70
5.2	Evaluasi Metode . . . . .	74
<b>6</b>	<b>KESIMPULAN DAN SARAN</b>	<b>77</b>
6.1	Kesimpulan Penelitian . . . . .	77
6.2	Saran Pengembangan . . . . .	77
	<b>DAFTAR REFERENSI</b>	<b>79</b>
	<b>A KODE PROGRAM</b>	<b>81</b>



## DAFTAR GAMBAR

1.1	Contoh pemodelan komunitas [3]	2
2.1	Karakteristik <i>big data</i> [4]	6
2.2	Teknologi <i>big data</i> dalam ekosistem Hadoop [4]	6
2.3	Contoh graf <i>complete</i> [1]	7
2.4	Contoh graf <i>cycle</i> [1]	7
2.5	Contoh graf <i>wheel</i> [1]	7
2.6	Contoh subgraf [1]	8
2.7	Komponen Apache Spark [5]	9
2.8	<i>Narrow vs Wide Transformation</i> [5]	10
2.9	Proses komputasi Spark [5]	11
2.10	Alur komputasi Spark [6]	11
2.11	Contoh visualisasi DAG	12
2.12	Contoh pola yang dicari <i>Motif Finding</i>	13
2.13	Potongan visualisasi DAG dari pencarian pola di atas	14
2.14	Tahapan eksekusi MF dalam <i>Optimized Plan Query</i>	14
2.15	Contoh bentuk komunitas SIC	15
2.16	Contoh bentuk komunitas SC	15
2.17	Contoh bentuk komunitas SCIC	16
2.18	Perbandingan hasil <i>stemming</i> dan lemmatisasi	16
2.19	Gambaran proses pemodelan topik	18
2.20	Contoh <i>document-term matrix</i>	19
2.21	Alur pemodelan topik dengan LSA	19
2.22	Proses dekomposisi <i>Document-Term Matrix</i> pada SVD	19
3.1	Contoh format data MEDLINE	22
3.2	Ilustrasi pemodelan graf	23
3.3	5 <i>subgraph</i> dengan simpul terbanyak	24
3.4	Visualisasi <i>subgraph</i> yang terbentuk	24
3.5	Topik yang terdapat pada <i>subgraph</i> raksasa	25
3.6	Topik yang terdapat pada <i>subgraph</i> lebih kecil	25
3.7	10 topik yang paling banyak mengalami <i>co-occurrence</i>	26
3.8	Detail data <i>node</i>	27
3.9	Detail data <i>edge</i>	27
3.10	Contoh data <i>node</i>	27
3.11	Contoh data <i>edge</i>	28
3.12	10 tanggal dengan rental terbanyak	28
3.13	Jumlah rental setiap bulan	29
3.14	Ringkasan cuaca Oktober-Desember 2014	30
3.15	10 rute sepeda terbanyak	31
3.16	10 rute sepeda terbanyak yang berawal dari atau berakhir di Townsend	31
3.17	Data sintetik eksperimen LSA	32
3.18	Contoh DTM <i>CountVectorizer</i>	33

3.19	Contoh DTM <i>TfidfVectorizer</i>	33
3.20	Contoh DF DT dengan bobot TF	34
3.21	Contoh DF DT dengan bobot TF-IDF	34
3.22	Contoh DF TT dengan bobot TF	34
3.23	Contoh DF TT dengan bobot TF-IDF	35
3.24	Hasil topik <i>latent</i> dengan bobot TF	35
3.25	Hasil topik <i>latent</i> dengan bobot TF-IDF	35
3.26	Detail data <i>node</i>	36
3.27	Detail data <i>edge</i>	36
3.28	Jumlah data <i>reply</i> dan <i>quote</i>	37
3.29	Perbandingan data <i>reply</i> dan <i>quote</i>	37
3.30	Skema akhir untuk <i>node</i> dan <i>edge</i>	38
3.31	Ilustrasi pemodelan graf	38
3.32	Visualisasi komunitas yang terbentuk	39
3.33	Contoh komunitas yang terbentuk	40
3.34	Akun <i>kompascom</i> yang menjadi salah satu pusat komunitas	41
3.35	Visualisasi SC yang terbentuk	41
3.36	Contoh data <i>tweet spam</i>	42
3.37	Contoh data <i>tweet spam(2)</i>	42
3.38	Contoh matriks <i>document-term</i>	43
3.39	Hasil nilai <i>singular</i> untuk $n = 2$	43
3.40	Hasil nilai <i>singular</i> untuk $n = 3$	43
3.41	Hasil nilai <i>singular</i> untuk $n = 4$	44
3.42	Hasil nilai <i>singular</i> untuk $n = 5$	44
3.43	Matriks distribusi dokumen per topik	44
3.44	Hasil nilai <i>singular</i> untuk kedua topik	45
4.1	Arsitektur sistem pengumpul data	48
4.2	Contoh data Twitter berformat JSON	48
4.3	Detail kolom dari salah satu sampel data	50
4.4	Contoh sampel data awal	51
4.5	Contoh sampel data <i>node</i>	51
4.6	Contoh sampel data <i>edge</i>	51
4.7	Contoh <i>retweet</i>	52
4.8	Contoh <i>quote</i>	52
4.9	Jumlah interaksi <i>reply</i> dan <i>quote</i>	53
4.10	Perbandingan jumlah interaksi <i>reply</i> dan <i>quote</i>	53
4.11	Tren jumlah <i>quote</i> dan <i>reply</i> data hidroponik	54
4.12	Tren jumlah <i>quote</i> dan <i>reply</i> data pemanasan global	55
4.13	Tren jumlah <i>quote</i> dan <i>reply</i> data kuliner indonesia	55
5.1	Detail data <i>node</i> hidroponik periode Mei-Juli	57
5.2	Detail data <i>edge</i> hidroponik periode Mei-Juli	58
5.3	Jumlah data unik interaksi <i>reply</i>	58
5.4	Jumlah data unik interaksi <i>quote</i>	58
5.5	Jumlah <i>reply</i> setelah menghapus <i>self loop</i>	59
5.6	Jumlah <i>quote</i> setelah menghapus <i>self loop</i>	59
5.7	Contoh postingan <i>thread</i>	59
5.8	Hasil <i>group by edge reply</i>	60
5.9	Hasil <i>group by edge quote</i>	60
5.10	Visualisasi komunitas yang terbentuk dari data <i>reply</i>	61
5.11	Visualisasi komunitas yang terbentuk dari data <i>quote</i>	61

5.12 Akun <i>tirtoid</i> . . . . .	62
5.13 Hasil deteksi <i>sc_1</i> pada data <i>quote</i> . . . . .	62
5.14 Hasil deteksi <i>sc_1</i> pada data <i>reply</i> . . . . .	63
5.15 Hasil deteksi <i>sc_2</i> pada data <i>quote</i> . . . . .	63
5.16 Hasil deteksi <i>sc_2</i> pada data <i>reply</i> . . . . .	63
5.17 Hasil deteksi <i>sic_1</i> pada data <i>quote</i> . . . . .	63
5.18 Hasil deteksi <i>sic_1</i> pada data <i>reply</i> . . . . .	63
5.19 Hasil deteksi <i>scic_1</i> pada data <i>quote</i> . . . . .	64
5.20 Hasil deteksi <i>scic_1</i> pada data <i>reply</i> . . . . .	64
5.21 Komunitas yang terdeteksi dengan anggotanya (dipisahkan <i>tab</i> ) pada data <i>reply</i> . . . . .	65
5.22 Sampel data hasil <i>join</i> . . . . .	65
5.23 Hasil pemodelan topik dari komunitas yang terdeteksi . . . . .	66
5.24 Detail data <i>node</i> kuliner Indonesia . . . . .	67
5.25 Detail data <i>edge</i> kuliner Indonesia . . . . .	67
5.26 Jumlah data unik <i>reply</i> data kuliner Indonesia . . . . .	67
5.27 Jumlah data unik <i>quote</i> data kuliner Indonesia . . . . .	67
5.28 Jumlah data unik <i>reply</i> data kuliner Indonesia tanpa <i>self loop</i> . . . . .	67
5.29 Jumlah data unik <i>quote</i> data kuliner Indonesia tanpa <i>self loop</i> . . . . .	68
5.30 Hasil <i>group by</i> data <i>reply</i> data kuliner Indonesia . . . . .	68
5.31 Hasil <i>group by</i> data <i>quote</i> data kuliner Indonesia . . . . .	68
5.32 Contoh komunitas yang terbentuk dari data <i>reply</i> . . . . .	69
5.33 Contoh komunitas yang terbentuk dari data <i>quote</i> . . . . .	69
5.34 <i>Tweet spam</i> yang diposting <i>pegipegi</i> dan mengandung kata "kuliner" . . . . .	70
5.35 Detail data <i>node</i> pemanasan global . . . . .	71
5.36 Detail data <i>edge</i> pemanasan global . . . . .	71
5.37 Jumlah data unik interaksi <i>reply</i> data pemanasan global . . . . .	71
5.38 Jumlah data unik interaksi <i>quote</i> data pemanasan global . . . . .	71
5.39 Jumlah data unik interaksi <i>reply</i> data pemanasan global tanpa <i>self loop</i> . . . . .	71
5.40 Jumlah data unik interaksi <i>quote</i> data pemanasan global tanpa <i>self loop</i> . . . . .	72
5.41 Hasil <i>group by</i> data <i>reply</i> pemanasan global . . . . .	72
5.42 Hasil <i>group by</i> data <i>quote</i> pemanasan global . . . . .	72
5.43 Visualisasi komunitas data <i>reply</i> pemanasan global . . . . .	73
5.44 Visualisasi komunitas data <i>quote</i> pemanasan global . . . . .	73
5.45 Contoh data sintetik <i>node</i> untuk pengujian . . . . .	74
5.46 Contoh data sintetik <i>edge</i> untuk pengujian . . . . .	74
5.47 Visualisasi SIC data sintetik . . . . .	75
5.48 Hasil pembersihan data teks . . . . .	75
5.49 Hasil pemodelan topik dengan jumlah 2 topik . . . . .	76

## DAFTAR KODE PROGRAM

3.1	Kode untuk membaca data XML . . . . .	22
3.2	Kode untuk <i>filtering</i> topik <i>major</i> . . . . .	22
3.3	Kode eksplorasi jumlah topik <i>major</i> . . . . .	22
3.4	Kode untuk <i>generate co-occurrence</i> . . . . .	23
3.5	Kode pembuatan RDD simpul dan sisi . . . . .	23
3.6	Kode analisis <i>Connected Components</i> . . . . .	24
3.7	Kode untuk analisis <i>subgraph</i> raksasa . . . . .	24
3.8	Kode untuk analisis <i>subgraph</i> kecil . . . . .	25
3.9	Kode <i>Degree Distribution</i> untuk melihat topik dengan <i>co-occurrence</i> terbanyak . . . . .	25
4.1	Kode untuk membaca data menggunakan <i>jsonlines</i> . . . . .	49
4.2	Kode untuk melakukan seleksi fitur . . . . .	49
4.3	Kode untuk memisahkan data <i>node</i> dan <i>edge</i> dan konversi ke Pandas <i>Dataframe</i> . . . . .	51
5.1	Kode <i>filtering</i> interaksi <i>reply</i> dan <i>quote</i> data hidroponik . . . . .	58
5.2	Kode pembuangan <i>self loop</i> data hidroponik . . . . .	59
5.3	Kode <i>group by edge</i> data hidroponik . . . . .	60
5.4	Kode konstruksi graf data hidroponik . . . . .	62
A.1	Eksperimen Data Sitasi MEDLINE . . . . .	81
A.2	Eksperimen Data Rute Sepeda . . . . .	82
A.3	Eksperimen Latent Semantic Analysis . . . . .	83
A.4	Eksperimen Data Twitter Covid Interaksi Quote . . . . .	84
A.5	Eksperimen Data Studi Kasus . . . . .	86

# BAB 1

## PENDAHULUAN

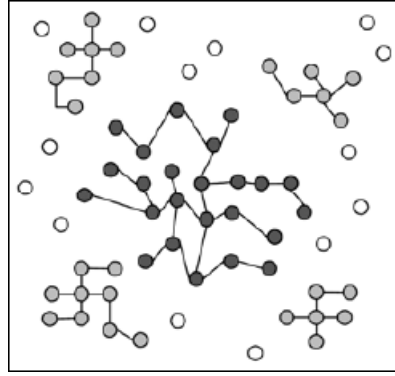
Pada bab ini akan dibahas tentang latar belakang penelitian yang mendeskripsikan gambaran besar permasalahan yang ada. Dilanjutkan dengan pembahasan rumusan masalah, tujuan penelitian, dan batasan penelitian. Pada tiga bagian tersebut, dijelaskan inti permasalahan, tujuan yang ingin dicapai sebagai solusi masalah, serta asumsi yang digunakan untuk membatasi ruang lingkup penelitian ini. Pembahasan tentang metodologi dan sistematika pembahasan menjadi penutup dari bab ini. Dua bagian tersebut menjelaskan tentang tahapan eksperimen dan penjabaran konten setiap bab pada buku ini.

### 1.1 Latar Belakang

Istilah *Big Data* merujuk pada sebuah himpunan data yang memiliki karakteristik khusus. Umumnya, karakteristik tersebut didefinisikan sebagai 3V, yaitu *Volume*, *Velocity*, dan *Variety*. *Volume* mendefinisikan ukuran data yang sangat besar, *Velocity* mendefinisikan penambahan data yang terjadi terus menerus dalam waktu singkat, sedangkan *Variety* mendefinisikan bentuk data yang sangat beragam [7]. Data dapat berupa data terstruktur seperti tabel basis data, data semi-terstruktur seperti XML dan JSON, dan data tidak terstruktur seperti data teks dari Twitter atau data citra dari satelit. Agar kumpulan data tersebut lebih bermakna, muncul suatu kebutuhan untuk menganalisisnya. Lalu, mencari informasi berharga dari kumpulan data tersebut, dan membuat keputusan/menarik kesimpulan berdasarkan hasil analisis data.

*Big data* yang datanya mendeskripsikan relasi antar objek, dapat diolah dan dimodelkan dalam bentuk *big graph*. *Big graph* adalah himpunan simpul dan sisi yang dibentuk dari *big data*. Sebagai contoh, data interaksi pengguna Twitter termasuk *big data* dapat dimodelkan dalam bentuk graf. Pemodelan graf menggunakan graf berarah (*directed graph*). Simpul pada graf memodelkan sebuah akun pengguna, sedangkan sisi pada graf memodelkan interaksi antar pengguna. Interaksi yang dimaksud dapat berupa *retweet*, *quote*, atau *reply*. Arah pada sisi graf menunjukkan interaksi sebuah akun terhadap akun lainnya. Sebagai contoh, seorang pengguna *re-tweet* atau membalas *tweet* pengguna lain. Setelah model graf terbentuk, nantinya graf dapat dianalisis menggunakan berbagai teknik. Saat ini, permasalahan yang sering diteliti pada analisis *big graph* adalah masalah deteksi komunitas dan analisis profil komunitas. Komunitas pada data Twitter, dapat didefinisikan sebagai kumpulan pengguna yang aktif berinteraksi pada periode tertentu, karena adanya ketertarikan untuk membicarakan suatu topik[8]. Inti permasalahan deteksi komunitas adalah mencari/menemukan kelompok pengguna yang saling berinteraksi. Jika data pengguna dimodelkan dalam graf, komunitas akan terlihat seperti graf yang padat (memiliki jumlah simpul dan sisi yang lebih banyak). Gambar 1.1 menunjukkan contoh visualisasi komunitas berupa graf yang lebih padat.

Analisis komunitas menjadi banyak diteliti karena melimpahnya data dari media sosial, sehingga ada potensi untuk menggali informasi berharga dari data tersebut. Berdasarkan hasil studi literatur dari Moertini dan Adithia, serta Needham dan Hodler [8, 9], ditemukan bahwa berbagai algoritma deteksi komunitas, seperti *Connected Component (CC)*, *Strongly Connected Component (SCC)*, dan *Triangle Count* telah berhasil diimplementasikan. Namun, sebagian besar untuk graf tidak berarah,



Gambar 1.1: Contoh pemodelan komunitas [3]

atau membutuhkan input jumlah komunitas. Kedua kendala tersebut menyebabkan implementasi dengan algoritma yang ada cukup sulit untuk kasus komunitas di Twitter, karena komunitas yang terbentuk sangat dinamis dan hanya berlaku pada periode tertentu saja. Bisa saja terbentuk komunitas baru, atau komunitas yang terbentuk pada suatu periode dapat menghilang di periode lain. Terbentuknya atau hilangnya komunitas dipengaruhi oleh interaksi yang dilakukan pengguna. Selain itu, terdapat algoritma yang komputasinya lambat, seperti CC dan SCC, sehingga kurang cocok jika diterapkan untuk *big data*.

Pada penelitian Moertini dan Adithia [8], berhasil diciptakan solusi untuk deteksi komunitas pada *directed big graph* dari data Twitter. Solusi yang diusulkan dengan menggunakan *Motif Finding*, yaitu teknik pengenalan pola (*pattern matching*) yang diimplementasikan pada *library* GraphFrames di Apache Spark. Dilakukan perbandingan deteksi komunitas menggunakan algoritma *Strongly Connected Component* (SCC) dan teknik *Motif Finding*. Secara waktu eksekusi, didapatkan bahwa teknik *Motif Finding* memperoleh hasil yang lebih cepat dibandingkan SCC. Namun, teknik *Motif Finding*, perlu mendefinisikan pola-pola komunitas (motif) yang ingin dideteksi.

Salah satu teknik analisis yang dapat diterapkan pada *big graph* data Twitter adalah analisis *Co-Occurrence Networks*. *Co-Occurrence Networks* bermakna kemunculan suatu kejadian dengan frekuensi tinggi dalam data yang bermodel graf [2]. Contoh *co-occurrence* adalah, kemunculan kata “vaksin” dalam kumpulan teks *tweet* tentang Covid yang dikirimkan pengguna. Analisis *Co-Occurrence Networks*, dapat dilakukan terhadap teks *tweet* sekumpulan pengguna yang sering berinteraksi (komunitas). Tujuannya untuk mencari tahu kata-kata yang sering dikirimkan kumpulan pengguna tersebut saat saling berinteraksi. Hasil analisisnya dapat diinterpretasi dan dimanfaatkan, seperti promosi produk untuk komunitas tertentu, mendeteksi topik dari sebuah komunitas yang aktif berinteraksi, atau menganalisis profil komunitas untuk mencari tahu topik apa yang dibicarakan anggota komunitas, sehingga kumpulan pengguna menjadi aktif berinteraksi.

Salah satu teknologi yang menunjang proses analisis *big graph* adalah *library* GraphX dan GraphFrames pada Apache Spark. Kedua *library* tersebut mengimplementasikan berbagai algoritma graf yang dijalankan secara paralel pada teknologi Apache Spark. Algoritma yang diimplementasikan seperti *PageRank*, *Connected Components*, dan *Breadth First Search*. Spark mampu melakukan komputasi terdistribusi yang menyebabkan proses komputasi berlangsung dengan cepat dan efisien. Oleh karena itu, cocok untuk digunakan pada komputasi yang dilakukan secara iteratif [5], seperti pada analisis *big graph*.

Skripsi ini berfokus pada deteksi komunitas dan analisis *Co-Occurrence* untuk analisis profil komunitas dari data Twitter. Data studi kasusnya berupa data interaksi pengguna Twitter, beserta *tweet* yang dikirim. Data berformat JSON, dikumpulkan secara *near real time* memanfaatkan Apache Kafka, dan disimpan pada *Hadoop Distributed File System* (HDFS). Penjelasan arsitektur sistem pengumpul data dapat diakses pada Gambar 4.1.

Data Twitter yang telah disimpan, nantinya akan dieksplorasi dan ditransformasi hingga dapat dimodelkan menjadi graf. Pembuatan graf memanfaatkan *library* GraphFrames pada Apache Spark.

Setelah graf terbentuk, akan dilakukan deteksi komunitas menggunakan teknik *Motif Finding* yang diimplementasikan pada GraphFrames. Komunitas yang terdeteksi akan dianalisis profilnya dengan cara analisis *tweet* yang sering dikirimkan di komunitas (*co-occur*). Sebelum dianalisis, data teks akan dibersihkan dan diproses dahulu menggunakan *library* Pandas dan PySastrawi di Python. Analisis Co-Occurrence dilakukan dengan teknik pemodelan topik *Latent Semantic Analysis* yang diimplementasikan pada *library* Scikit-Learn di Python.

Perangkat lunak yang dihasilkan terdiri dari tiga bagian, yaitu program untuk melakukan eksplorasi dan penyiapan data, program untuk pembuatan graf dan deteksi komunitas, dan program untuk pemrosesan teks dan pemodelan topik. Program pertama dan kedua diimplementasikan dalam bahasa pemrograman Scala dan dijalankan dengan basis *command line interface*. Sedangkan program ketiga diimplementasikan dalam bahasa pemrograman Python dan dijalankan pada Google Colaboratory, yaitu sebuah teknologi berbasis *cloud* untuk mengeksekusi kode Python.

## 1.2 Rumusan Masalah

Adapun rumusan masalah dari latar belakang yang telah dipaparkan adalah:

1. Bagaimana cara implementasi deteksi komunitas dan analisis *Co-Occurrence Networks* dari data interaksi pengguna Twitter?
2. Bagaimana cara eksplorasi dan menyiapkan data interaksi pengguna Twitter hingga dapat dimodelkan dalam graf dan dapat dianalisis *co-occurrence*-nya?
3. Bagaimana cara analisis, dan interpretasi hasil analisis *big data* studi kasus?
4. Bagaimana cara membangun perangkat lunak yang mampu melakukan penyiapan dan analisis data studi kasus?

## 1.3 Tujuan

Adapun tujuan penelitian dari rumusan masalah yang telah dipaparkan adalah:

1. Mempelajari konsep deteksi komunitas dan analisis *Co-Occurrence Networks*;
2. Melakukan eksplorasi dan penyiapan *big data* studi kasus, merancang pemodelan grafnya, dan merumuskan aspek yang dapat dianalisis *co-occurrence*-nya;
3. Melakukan analisis *big data* studi kasus dan menginterpretasikan hasil analisisnya untuk menemukan informasi berharga;
4. Membangun perangkat lunak yang mampu melakukan penyiapan dan analisis *big data* studi kasus.

## 1.4 Batasan Penelitian

Batasan untuk penelitian ini adalah sebagai berikut:

1. Teks *tweet* yang dianalisis interaksinya, hanya berasal dari domain tertentu dan dalam Bahasa Indonesia;
2. Teks *tweet* yang dianalisis tidak dilakukan pemrosesan tambahan, seperti *n-gram*, lemmatisasi, *named entity recognition*, dan *parsing*, karena keterbatasan *library* untuk pengolahan teks Bahasa Indonesia.
3. Teks *tweet* yang *spam* tidak ditangani.
4. Interaksi antar pengguna yang dianalisis hanya *quote* atau *reply* saja.

## 1.5 Metodologi

Metodologi yang digunakan dalam pembuatan skripsi ini adalah sebagai berikut:

1. Melakukan eksplorasi dan studi literatur tentang Apache Spark;

2. Melakukan eksplorasi dan studi literatur tentang bahasa pemrograman Scala;
3. Melakukan eksplorasi dan studi literatur tentang *library* GraphX dan GraphFrames;
4. Melakukan studi literatur tentang konsep graf dan deteksi komunitas;
5. Melakukan eksplorasi dan studi literatur tentang konsep *Co-Occurrence Networks*;
6. Melakukan eksplorasi dan studi literatur tentang konsep pemodelan topik, *text pre-processing*, dan fungsi-fungsi terkait pada *library* Scikit-Learn dan PySastrawi;
7. Mengumpulkan *big data* studi kasus yang dapat ditransformasi menjadi *big graph* dan berpotensi untuk analisis *Co-Occurrence Networks*;
8. Melakukan eksplorasi dan penyiapan *big data* studi kasus;
9. Membuat perangkat lunak yang mampu melakukan penyiapan data, analisis data, dan menampilkan hasil analisisnya;
10. Melakukan analisis terhadap hasil yang didapatkan, merumuskan interpretasi hasil analisis, dan memaparkan hasilnya dengan menarik.

## 1.6 Sistematika Pembahasan

Sistematika penulisan skripsi ini adalah sebagai berikut.

1. Bab Pendahuluan
2. Bab Landasan Teori
3. Bab Studi Eksplorasi Apache Spark, GraphFrames, dan Scikit-Learn
4. Bab Pengumpulan, Eksplorasi, dan Penyiapan Data
5. Bab Analisis Data dan Evaluasi
6. Bab Kesimpulan dan Saran

Bab Pendahuluan membahas tentang pengantar skripsi yang berisi gambaran besar permasalahan yang dibahas, batasan ruang lingkup penelitian, serta tahapan eksperimen yang dilakukan. Bab Landasan Teori berisi pembahasan hasil studi literatur yang mencakup penjelasan berbagai konsep yang digunakan. Bab Studi Eksplorasi Teknologi membahas tentang berbagai eksperimen untuk mempelajari konsep dan teknologi yang digunakan pada skripsi ini. Pada Bab tiga juga dijelaskan tentang contoh penerapan metodologi penelitian (eksplorasi dan penyiapan data, pemodelan graf, deteksi komunitas dengan *Motif Finding*, tahap pra-pemrosesan data teks, dan pemodelan topik dengan *Latent Semantic Analysis*). Bab Pengumpulan, Eksplorasi, dan Penyiapan Data membahas tentang penerapan metodologi penelitian tahap pertama (pengumpulan, eksplorasi, penyiapan data, dan pemodelan graf) pada data studi kasus (detail data studi kasus dijelaskan pada bab terkait). Bab Analisis Data dan Evaluasi membahas tentang penerapan metodologi penelitian tahap selanjutnya (pembuatan graf, deteksi komunitas, pemodelan topik, dan interpretasi hasil analisis). Buku ini diakhiri dengan Bab Kesimpulan dan Saran.