

SISTEM OLAHAN TEKS DENGAN MENGGUNAKAN SOALAN TERBUKA

ERLINDA LUKE JERRY

Projek ini merupakan salah satu keperluan untuk
Ijazah Sarjana Muda Sains dengan Kepujian
Sains Kognitif

Fakulti Sains Kognitif dan Pembangunan Sumber Manusia
UNIVERSITI MALAYSIA SARAWAK
2004

PENGHARGAAN

Pertama sekali saya ingin mengucapkan rasa bersyukur dan penghargaan kepada ahli keluarga saya kerana telah memberikan sokongan penuh kepada saya sepanjang saya menjalani Projek Tahun Akhir ini.

Dengan kesempatan ini, saya juga ingin mengucapkan ribuan terima kasih kepada Encik Syafiq Fikri Abdullah, selaku Penyelia Projek Tahun Akhir saya, yang telah banyak membantu saya bagi membolehkan saya menghasilkan laporan akhir Projek Tahun Akhir ini. Beliau juga telah banyak membantu di dalam memberikan garis panduan di dalam pembangunan sistem yang saya bangunkan bagi Projek Tahun Akhir saya ini.

Setinggi-tinggi penghargaan dan ucapan terima kasih yang tidak terhingga ini juga saya ucapkan kepada rakan-rakan saya yang banyak membantu di dalam memberikan sokongan dan nasihat sepanjang saya menjalani Projek Tahun Akhir ini. Tidak lupa juga ucapan terima kasih ini saya tujukan kepada responden saya di atas kerjasama mereka di dalam menjawab soalan terbuka bagi kajian ini.

Oleh itu, sekali lagi saya ingin mengucapkan ribuan terima kasih kepada mereka yang telah banyak membantu saya sepanjang saya menjalankan Projek Tahun Akhir ini sama ada secara langsung atau tidak langsung. Kejayaan saya dalam menjalankan Projek Tahun Akhir ini dan dengan terhasilnya laporan akhir Projek Tahun Akhir ini adalah atas dorongan serta sokongan daripada mereka terhadap saya.

JADUAL KANDUNGAN

| | |
|----------------------------------------------|------|
| Penghargaan | iii |
| Jadual kandungan | iv |
| Senarai Jadual | vi |
| Senarai Rajah | viii |
| Abstrak | x |
| <i>Abstract</i> | xi |
| | |
| 1. Pendahuluan | |
| 1.1 Pengenalan | 1 |
| 1.2 Latar Belakang Kajian | 2 |
| 1.3 Pernyataan Masalah | 3 |
| 1.4 Objektif Kajian | |
| 1.4.1 Objektif Umum | 4 |
| 1.4.2 Objektif Khusus | 4 |
| 1.5 Skop Kajian | 5 |
| 1.6 Kepentingan Kajian | 5 |
| 1.7 Definisi istilah | |
| 1.7.1 Olahan teks | 6 |
| 1.7.2 Olahan data | 6 |
| 1.7.3 Soalan terbuka | 6 |
| 1.7.4 "Classification Rule" | 6 |
| 1.7.5 "Association Rule" | 7 |
| 1.8 Limitasi Kajian | 7 |
| 1.9 Sinopsis Kajian | 7 |
| | |
| 2. Sorotan Kajian Lepas | |
| 2.1 Pengenalan | 9 |
| 2.2 Olahan Data | 9 |
| 2.3 Olahan Teks | 11 |
| 2.4 Soalan Terbuka | 13 |
| 2.5 "Classification Rule" | 14 |
| 2.6 "Association Rule" | 16 |
| | |
| 3. Metodologi Sistem | |
| 3.1 Pengenalan | 19 |
| 3.2 Metodologi | 19 |
| 3.2.1 Pernyataan masalah dan objektif kajian | 21 |
| 3.2.2 Sorotan kajian lepas | 21 |
| 3.2.3 Mengenalpasti keperluan sistem | 21 |
| 3.2.4 Pembentukan soalan terbuka | 22 |
| 3.2.5 Algoritma dan Rekabentuk sistem | 22 |
| 3.2.6 Implementasi Sistem | 23 |
| 3.2.7 Pengujian sistem | 23 |
| 3.2.8 Keputusan dan analisis | 24 |

| | | |
|------|---------------------------------------------------------------|----|
| 4. | Keputusan Kajian | |
| 4.1 | Pengenalan | 25 |
| 4.2 | Proses-Proses Pengiraan Skor Bagi Katakunci | 25 |
| | 4.2.1 Contoh Pengiraan Untuk “Classification Rule” | 28 |
| | 4.2.2 Contoh Pengiraan Untuk “Association Rule” | 32 |
| 4.3 | “Classification rule” Bagi Motorola | 33 |
| 4.4 | “Classification rule” Bagi Nokia | 34 |
| 4.5 | “Classification rule” Bagi Samsung | 35 |
| 4.6 | “Classification rule” Bagi Siemen | 36 |
| 4.7 | “Association rule” Bagi Motorola | 37 |
| 4.8 | “Association rule” Bagi Nokia | 38 |
| 4.9 | “Association rule” Bagi Samsung | 38 |
| 4.10 | “Association rule” Bagi Siemen | 39 |
| 4.11 | Graf “Classification Rule” Bagi Motorola | 40 |
| 4.12 | Graf “Classification Rule” Bagi Nokia | 41 |
| 4.13 | Graf “Classification Rule” Bagi Samsung | 42 |
| 4.14 | Graf “Classification Rule” Bagi Siemen | 43 |
| 4.15 | Graf “Association Rule” Bagi Motorola | 44 |
| 4.16 | Graf “Association Rule” Bagi Nokia | 45 |
| 4.17 | Graf “Association Rule” Bagi Samsung | 46 |
| 4.18 | Graf “Association Rule” Bagi Siemen | 47 |
| 4.19 | Histogram Bagi Motorola | 48 |
| 4.20 | Histogram Bagi Nokia | 48 |
| 4.21 | Histogram Bagi Samsung | 49 |
| 4.22 | Histogram Bagi Siemen | 50 |
| 4.21 | Data “Co-occurrence” Untuk Katakunci Dan Jenis Telefon Bimbit | 50 |
| 5. | Kesimpulan | |
| 5.1 | Pengenalan | 52 |
| 5.2 | Kesimpulan Kajian | 52 |
| 5.3 | Kebaikan Kajian | 54 |
| 5.4 | Cadangan kerja-kerja pada masa akan datang | |
| | 5.4.1 Jawapan Kepada Soalan Terbuka Secara “On-line” | 54 |
| | 5.4.2 Menambahkan bilangan responden | 55 |
| | 5.4.3 Jadual Keputusan Penilaian | 55 |
| | 5.4.4 “Positioning Map” | 55 |
| 6. | Rujukan | 56 |
| 7. | Lampiran A – Contoh Soalan Terbuka | 58 |
| 8. | Lampiran B – Implementasi Sistem | 60 |

SENARAI JADUAL

| | |
|------------------------------------------------------------------------------|----|
| Jadual 2.1 Aplikasi-aplikasi olahan data di dalam dunia perniagaan | 10 |
| Jadual 2.2 Keputusan Penilaian | 13 |
| Jadual 2.3 Contoh data soal selidik | 14 |
| Jadual 2.4 Histogram untuk imej kereta jenama A | 15 |
| Jadual 2.5 “Classification rule” bagi kereta jenis A | 15 |
| Jadual 2.6 Lima parameter bagi “Association Rules” | 16 |
| Jadual 4.1 Keputusan “Classification rule” bagi Motorola | 33 |
| Jadual 4.2 Keputusan “Classification rule” bagi Nokia | 34 |
| Jadual 4.3 Keputusan “Classification rule” bagi Samsung | 35 |
| Jadual 4.4 Keputusan “Classification rule” bagi Siemen | 36 |
| Jadual 4.5 Keputusan “Association rule” bagi Motorola | 37 |
| Jadual 4.6 Keputusan “Association rule” bagi Nokia | 38 |
| Jadual 4.7 Keputusan “Association rule” bagi Samsung | 38 |
| Jadual 4.8 Keputusan “Association rule” bagi Siemen | 39 |
| Jadual 4.9 Histogram Bagi Motorola | 48 |

| | |
|-------------------------------------------------------------------------------------|----|
| Jadual 4.10 Histogram Bagi Nokia | 48 |
| Jadual 4.11 Histogram Bagi Samsung | 49 |
| Jadual 4.12 Histogram Bagi Siemen | 50 |
| Jadual 4.13 Data “Co-occurrence” untuk katakunci dan jenis telefon bimbit | 50 |

Demo (Visit <http://www.pdfsplitmerger.com>)

SENARAI RAJAH

| | | |
|------------------|-----------------------------------------------------------------------|----|
| Rajah 2.1 | Perkaitan antara “Relevant documents” dan “Retrieved documents” | 12 |
| Rajah 2.2 | Contoh bagi Faktor “Confidence” dan Faktor “Support” | 17 |
| Rajah 3.1 | Lapan fasa untuk sistem olahan teks dengan menggunakan soalan terbuka | 20 |
| Rajah 4.1 | Peraturan Pemilihan dengan menggunakan SC | 27 |
| Rajah 4.2 | Graf “Classification rule” bagi Jadual 4.1 | 40 |
| Rajah 4.3 | Graf “Classification rule” bagi Jadual 4.2 | 41 |
| Rajah 4.4 | Graf “Classification rule” bagi Jadual 4.3 | 42 |
| Rajah 4.5 | Graf “Classification rule” bagi Jadual 4.4 | 43 |
| Rajah 4.6 | Graf “Classification rule” bagi Jadual 4.5 | 44 |
| Rajah 4.7 | Graf “Classification rule” bagi Jadual 4.6 | 45 |
| Rajah 4.8 | Graf “Classification rule” bagi Jadual 4.7 | 46 |
| Rajah 4.9 | Graf “Classification rule” bagi Jadual 4.8 | 47 |
| Rajah B.1 | Antaramuka “frmIntro.frm” | 60 |
| Rajah B.2 | Mesej Ralat 1 | 61 |
| Rajah B.3 | Antaramuka “frmMenu.frm” | 61 |

| | |
|--------------------------------------------------------|----|
| Rajah B.4 Mesej Ralat 2 | 62 |
| Rajah B.5 Antaramuka “frmLoading.frm” | 62 |
| Rajah B.6 Antaramuka “frmmain.frm” | 63 |
| Rajah B.7 Antaramuka “frmClassification.frm” | 64 |
| Rajah B.8 Antaramuka “frmAssociation.frm” | 66 |
| Rajah B.9 Antaramuka “frmCoocurrence.frm” | 67 |

Demo (Visit <http://www.pdfsplitmerger.com>)

ABSTRAK

Kajian ini bertujuan untuk membuat kajian mengenai olahan teks dengan menggunakan soalan terbuka. Proses menganalisis soalan terbuka adalah merupakan tugas yang agak sukar kerana jawapan kepada soalan terbuka adalah tidak terhad yang mana ia membenarkan responden untuk memberikan jawapan mereka secara bebas. Untuk mengatasi masalah ini, satu sistem olahan teks telah dibangunkan iaitu Sistem Olahan Teks Dengan Menggunakan Soalan Terbuka. Input untuk sistem ini diperolehi daripada jawapan kepada soalan terbuka yang melibatkan seramai 40 orang responden yang terdiri daripada pelajar-pelajar di Unimas. Terdapat empat jenama telefon bimbit yang dipilih untuk kajian ini. Ia adalah Motorola, Nokia, Samsung dan Siemen. Responden dikehendaki untuk memberikan jawapan mengenai ciri-ciri telefon bimbit yang mereka miliki. Di dalam kajian ini, terdapat dua hukum yang digunakan untuk menganalisis teks bagi soalan terbuka iaitu "classification rule" dan "association rule". Sistem ini dibangunkan dengan menggunakan Microsoft Visual Basic 6. Keputusan daripada kajian ini dapat membantu di dalam memahami proses-proses yang berlaku di dalam menganalisis soalan terbuka dengan menggunakan "classification rule" dan "association rule".

Demo (Visit <http://www.pdfsplitmerg.com>)

ABSTRACT

The purpose of this study is to do the research about text mining with open question. The process of analyzing open question is a difficult task because the answer to the open question is not fixed whereby it allows the respondents to answer it freely with their own answer. To overcome this problem, one system has been developed namely the Text Mining With Open Question System. The input for this system is taken from the answer for the open question that involves 40 respondents of Unimas's students. There are four types of handphone being chosen for this study. They are Motorola, Nokia, Samsung and Siemen. The respondents are required to give their answers about the characteristics of their handphone. In this study, there are two rules being used to analyze the text for open answer namely classification rule and association rule. The system is developed by using Microsoft Visual Basic 6. The result of this study will help the understanding of the process involved in analyzing the open question with the use of classification rule and association rule.

Demo (Visit <http://www.pdfsplitmerger.com>)

BAB 1

PENDAHULUAN

1.1 Pengenalan

Olahan data bukanlah merupakan suatu perkataan baru lagi di dalam dunia teknologi maklumat pada masa kini. Ia telah dipelajari di seluruh dunia dan dipelopori terutama sekali oleh mereka yang mempelajari atau mengkaji berkaitan dengan bidang tersebut. Selain olahan data, konsep olahan teks juga telah menjadi bertambah popular sebagai alat bagi pengurusan maklumat yang dikatakan mampu mendedahkan struktur maklumat yang boleh membantu di dalam proses mendapatkan keputusan yang pasti. Olahan teks melihat kepada corak-corak di dalam bahasa teks biasa. Ia didefinisikan sebagai proses-proses menganalisa teks untuk mengambil maklumat daripadanya untuk tujuan tertentu.

Terdapat pelbagai jenis kajian yang telah dibuat yang melibatkan olahan data dan olahan teks. Sebelum ini terdapat kajian yang dibuat terhadap soalan tertutup. Seterusnya, kajian telah dibuat ke atas soalan terbuka dan kajian ini adalah lebih sukar berbanding kajian yang dibuat bagi soalan tertutup. Pembangunan sistem bagi kajian ini juga bukanlah sesuatu yang mudah. Ini adalah kerana lebih mudah bagi sistem untuk mengecam teks bagi soalan tertutup berbanding teks bagi soalan terbuka. Kajian dan sistem yang dibangunkan oleh Yamanishi dan Li (2001), turut menjadi perhatian apabila mereka berjaya mencipta sistem olahan teks yang bertajuk "Mining Open Answer in Questionnaire Data".

Thesis ini bertujuan untuk membuat kajian dan seterusnya membangunkan sebuah sistem olahan teks bagi mengecam teks yang melibatkan soalan terbuka dengan menggunakan “classification rule” dan “association rule”.

1.2 Latar Belakang Kajian

Di dalam dunia mengejar era informasi bermaklumat, terdapat pelbagai kajian telah dilakukan untuk mendapatkan jawapan bagi permasalahan yang wujud berkaitan dengan dunia IT. Saban hari semakin ramai pakar-pakar komputer tampil dengan ciptaan-ciptaan terbaru mereka dan ini mewujudkan persaingan yang hebat di antara mereka. Perkataan teks bukanlah merupakan suatu perkataan yang baru bagi semua orang. Pelbagai kajian dan ciptaan telah dihasilkan untuk mengecam teks.

Di dalam pembangunan sistem olahan teks ini, beberapa soalan terbuka yang berkaitan dengan ciri-ciri yang disukai pada sesebuah telefon bimbit responden telah dibuat. Ini bertujuan untuk mendapatkan input bagi sistem yang akan dibangunkan. Melalui jawapan yang diperolehi daripada soalan terbuka tersebut, ia akan dijadikan sebagai input untuk membantu di dalam pengujian sistem dan seterusnya untuk mengenalpasti sama ada sistem yang dibangunkan adalah berjaya atau tidak.

Kebaikan soalan terbuka adalah ia membolehkan responden memberikan jawapan mengikut pendapat serta pandangan mereka sendiri tanpa terikat kepada jawapan-jawapan yang telah ditetapkan seperti yang terdapat pada soalan tertutup. Responden bukan saja boleh memberikan jawapan dalam bentuk teks, malah jawapan yang diberikan boleh juga dalam bentuk gambarajah mengikut pemahaman responden terhadap soalan terbuka yang diberikan. Ini membuka ruang

kepada responden untuk memberikan idea-idea mereka yang tersendiri berkaitan dengan jawapan kepada soalan terbuka yang diberikan dan ini sekaligus boleh membantu di dalam mendapatkan jawapan yang bernas.

Di samping kebaikan, terdapat juga keburukan soalan terbuka. Masalah yang sering timbul ialah apabila jawapan yang diberikan adalah tidak menepati kehendak soalan. Kadangkadang terdapat juga responden yang tidak memberikan kerjasama sepenuhnya di dalam menjawab soalan terbuka yang diberikan. Keadaan ini menyukarkan pengkaji untuk mendapatkan jawapan yang bersesuaian dan jawapan yang menepati kehendak pengkaji.

1.3 Pernyataan Masalah

Terdapat beberapa masalah yang timbul yang mendorong pengkaji untuk membuat kajian berkaitan olahan teks. Masalah-masalah yang dimaksudkan adalah seperti berikut:

- a) Kekaburan di dalam pengertian sebenar olahan data, olahan teks, “classification rule”, “association rule” dan soalan terbuka.
- b) Masalah di dalam proses pengiraan yang melibatkan “classification rule” dan “association rule”.
- c) Masalah di dalam pembangunan sistem berkaitan dengan pengecaman teks yang berkaitan dengan soalan terbuka.

1.4 Objektif Kajian

Objektif kajian ini terbahagi kepada dua bahagian, iaitu objektif umum dan objektif khusus.

1.4.1 Objektif Umum

Kajian ini bertujuan untuk mengkaji olahan teks (“text mining”) dengan menggunakan soalan terbuka.

1.4.2 Objektif Khusus

Terdapat 4 objektif khusus dalam kajian ini iaitu :

- Mempelajari serta memahami konsep olahan teks yang melibatkan soalan terbuka.
- Mengaplikasikan pengetahuan mengenai olahan teks dengan soalan terbuka di dalam pembangunan sistem olahan teks dengan menggunakan perisian Visual Basic.
- Membangunkan sebuah sistem olahan teks yang melibatkan soalan terbuka berkaitan dengan ciri-ciri yang disukai pada telefon bimbit yang merangkumi empat jenis telefon bimbit iaitu Motorola, Nokia, Siemen dan Samsung.
- Mencari “classification rule” dan “association rule” bagi beberapa jenis telefon bimbit.

1.5 Skop Kajian

Di dalam kajian ini, terdapat beberapa skop yang diharapkan dapat membantu dalam memahami dengan lebih berkaitan olahan teks. Berikut adalah skop-skop kajian yang dibuat :

- Pengertian olahan data, olahan teks, “classification rule”, “association rule” dan soalan terbuka.
- Pembentukan soalan terbuka.
- Merekabentuk antaramuka untuk pra-pemprosesan.
- Pembangunan sistem olahan teks dengan menggunakan perisian Microsoft Visual Basic 6.

1.6 Kepentingan Kajian

Dengan terhasilnya sistem olahan teks ini dapat membantu di dalam pengecaman teks yang berkaitan dengan soalan terbuka. Sistem ini dibuat memandangkan untuk pengecaman teks bagi soalan terbuka adalah agak sukar jika dibandingkan dengan pengecaman tek yang melibatkan soalan tertutup. Melalui kajian ini juga akan lebih membantu di dalam memahami proses-proses yang berlaku bagi kedua-dua hukum yang terlibat di dalam olahan teks yang melibatkan soalan terbuka iaitu “classification rule” dan “association rule”. Sistem yang dibangunkan pula dapat membantu untuk lebih memahami lagi proses-proses yang berlaku di dalam olahan teks bagi soalan terbuka.

1.7 Definisi Istilah

1.7.1 Olahan Teks

Menurut Witten dan Frank (2000), olahan teks adalah melihat kepada pengecaman corak dalam bentuk teks. Di dalam konteks ini, kita tidak perlu untuk memahami teks bagi mendapatkan maklumat-maklumat yang berguna yang terkandung di dalam teks berkenaan.

1.7.2 Olahan Data

Menurut Witten dan Frank (2000), olahan data adalah melihat kepada pengecaman corak dalam bentuk data. Ia adalah berkaitan dengan penyelesaian masalah dengan menganalisa data yang terdapat di dalam pangkalan data.

1.7.3 Soalan Terbuka

Soalan terbuka adalah merupakan soalan yang meminta jawapan yang panjang dan jawapan adalah tidak terdapat di dalam senarai jawapan. Jawapan boleh mengandungi teks, numerik atau gambarajah.

1.7.4 "Classification Rule"

Ia mengelaskan kes-kes kepada beberapa kelas. Menurut Yi dan Yamanishi (2001), "classification rules" mengandungi satu rangkaian atau siri mengikut tertib atau turutan peraturan "IF-THEN-ELSE" untuk menentukan atau menetapkan jawapan terbuka kepada sasaran.

1.7.5 “Association Rule”

Ia mengecam satu kombinasi bagi nilai bagi sifat atau perkara-perkara yang berlaku secara serentak dengan frekuensi yang tinggi berbanding yang sepatutnya jika nilai bagi sifat atau perkara-perkara adalah tidak bergantung di antara satu sama lain.

1.8 Limitasi Kajian

Limitasi-limitasi yang dikenalpasti dalam kajian ini ialah kajian ini hanya melibatkan golongan pelajar di Unimas sahaja bagi menjawab soalan terbuka yang dibuat bagi mendapatkan input bagi sistem yang dibangunkan. Kakitangan-kakitangan akademik dan kakitangan-kakitangan bukan akademik tidak termasuk di dalam menjawab soalan terbuka ini. Kesahihan jawapan kepada soalan terbuka yang diedarkan kepada responden bagi kajian ini adalah bergantung kepada kejujuran responden di dalam memberikan jawapan bagi soalan terbuka yang diberikan. Akhir sekali, bilangan responden atau pelajar yang terlibat di dalam menjawab soalan terbuka ini adalah seramai 40 orang.

1.9 Sinopsis Kajian

Bahagian ini mengandungi garis kasar mengenai laporan bagi kajian ini yang terdiri daripada lima bab. Setiap bab akan menggariskan skop yang berbeza yang terlibat di dalam kajian ini.

Bab 1 akan menyentuh mengenai pengenalan kepada kajian. Di dalam bab ini juga akan membantu pembaca untuk memahami berkaitan pengertian olahan data, olahan teks,

“classification rule” dan “association rule”, latar belakang kajian, pernyataan masalah, objektif kajian, skop kajian, metodologi kajian, kepentingan kajian, definisi istilah dan limitasi kajian.

Bab 2 pula mengandungi sorotan kajian lepas yang berkaitan dengan kajian ini.

Bab 3 akan menyentuh mengenai metodologi yang digunakan di dalam pembangunan sistem olahan teks. Ini akan membantu pembaca di dalam memahami dengan lebih lagi berkaitan dengan proses-proses yang terlibat di dalam kajian ini. Bab ini menyentuh tentang fasa-fasa yang terlibat di dalam kajian ini.

Bab 4 pula mengandungi keputusan kajian yang telah dibuat. Ia menyentuh mengenai hasil-hasil yang diperolehi daripada sistem yang telah dibangunkan di dalam kajian ini.

Bab 5 mengandungi kesimpulan keseluruhan bagi kajian ini dan cadangan-cadangan untuk memperbaiki sistem ini pada masa akan datang.

BAB 2

SOROTAN KAJIAN LEPAS

2.1 Pengenalan

Bab ini akan membicarakan tentang olahan teks, olahan data, soalan terbuka, “classification rule” dan “association rule” berdasarkan kajian lepas yang telah dibuat yang berkaitan dengan kajian ini.

2.2 Olahan Data

Olahan data lebih melihat kepada corak yang berbentuk data. Nama lain bagi olahan data ialah “Knowledge Discovery in Databases - KDD”. Ia menyentuh tentang penyelesaian masalah dengan menganalisa data yang terdapat di dalam pangkalan data. Menurut Witten dan Frank (2000), olahan data adalah merupakan proses mendapatkan corak yang terdapat di dalam data. Proses yang terlibat mestilah secara automatik atau separa automatik. Selain itu, corak yang dicari mestilah corak yang berguna yang boleh membawa kepada beberapa kebaikan.

Menurut Jessop (2001), olahan data ditakrifkan sebagai satu proses yang dikawal oleh komputer dengan mendapatkan arah, corak dan perhubungan yang sistematik dari stor data yang tidak dikenalpasti sebelumnya. Menurut Holshemier dan Siebes (1994) pula, olahan data adalah merupakan pencarian perhubungan dan corak-corak yang global yang wujud di dalam pangkalan data yang besar tetapi tersembunyi di antara bilangan data yang sangat luas.

Penggunaan olahan data adalah sangat meluas. Pelbagai aplikasi bagi olahan data telah dikenalpasti. Menurut Ponniah dan Wiley (2001), terdapat pelbagai aplikasi yang membawa

kepada keuntungan daripada olahan data. Jadual 2.1 menunjukkan beberapa aplikasi olahan data di dalam dunia perniagaan yang diambil daripada sumber Ponniah dan Wiley (2001).

| | |
|--------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| “Customer Segmentation” | Dunia perniagaan menggunakan olahan data untuk memahami pelanggan mereka. Algoritma Pengesanan Gugusan mendapatkan gugusan pelanggan dengan berkongsi ciri-ciri yang sama. |
| “Market Basket Analysis” | Ia merupakan aplikasi yang berguna untuk perniagaan runcit. Algoritma Analisis Hubungan membuka pertalian antara produk-produk yang dibeli pada masa yang sama. Perniagaan-perniagaan lain seperti “Upscale Auction Houses” menggunakan algoritma ini untuk mencari pelanggan-pelanggan untuk mengenalpasti kepada siapa mereka boleh menjual barang-barang yang bermutu tinggi. |
| “Risk Management” | Syarikat-syarikat insuran dan gadai janji (mortgage) menggunakan olahan data untuk membuka risiko-risiko yang berkaitan dengan pelanggan-pelanggan yang berpotensi. |
| “Fraud Detection” | Syarikat-syarikat Kad Kredit menggunakan olahan data untuk mendapatkan corak perbelanjaan pelanggan-pelanggan yang tidak normal. Corak-corak seperti ini boleh mendedahkan penggunaan palsu bagi kad kredit. |
| “Delinquency Tracking” | Syarikat-syarikat pinjaman menggunakan teknologi mengesan jejak pelanggan-pelanggan yang gagal memenuhi kewajiban untuk pembayaran balik pinjaman mereka. |
| “Demand Prediction” | Jualan runcit dan perniagaan-perniagaan lain menggunakan olahan data untuk menyesuaikan tuntutan dan membekalkan haluan / arah tuntutan untuk produk-produk yang khusus. |

Jadual 2.1 : Aplikasi-aplikasi olahan data di dalam dunia perniagaan

Penggunaan olahan data banyak membantu di dalam mengenalpasti hala tuju bagi sesebuah syarikat. Di antaranya adalah untuk mengenalpasti sebab-sebab yang mendorong pelanggan untuk membeli sesuatu produk. Olahan data juga membantu di dalam memberikan idea-idea untuk cara pemasaran secara terus. Selain itu, ia juga membantu di dalam mengenalpasti cara melatih pekerja-pekerja.

Menurut Jessop (2001), terdapat tiga langkah-langkah untuk olahan data. Langkah pertama ialah pengumpulan data, iaitu data dikumpul daripada pangkalan data yang berbeza kepada satu gudang. Proses ini selalunya dirujuk sebagai “data warehousing”. Langkah kedua pula ialah pembentukan model iaitu termasuk pembelajaran dan latihan. Model ini digunakan untuk membuat ramalan. Langkah yang ketiga pula ialah pengesahan model yang mana ia adalah merupakan proses pengujian terhadap model untuk memastikan ketepatannya.

2.3 Olahan Teks

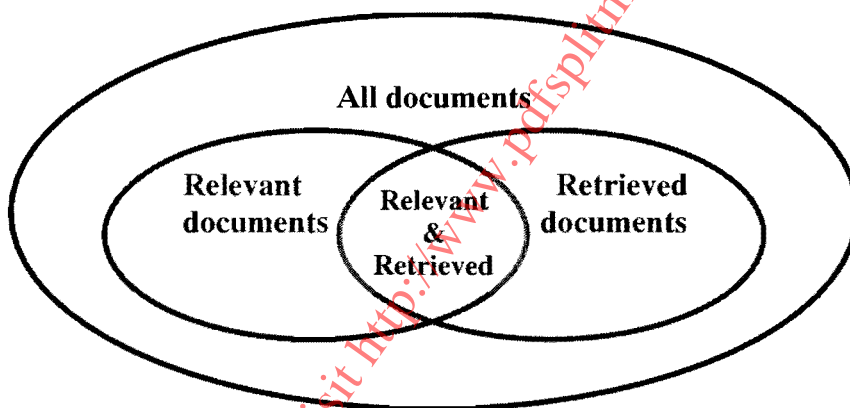
Olahan teks lebih melihat kepada corak yang berbentuk teks. Ia merupakan proses menganalisa teks untuk mendapatkan maklumat yang berguna bagi sesuatu tujuan. Teks adalah merupakan sesuatu yang tidak berstruktur dan ia juga tidak mempunyai bentuk. Tanpa kita sedari, teks telah dijadikan sebagai alat yang membantu di dalam proses pertukaran maklumat formal. Olahan teks adalah mudah kerana kita tidak perlu memahami teks. Secara umumnya, olahan teks mengandungi analisis tentang dokumen-dokumen teks dengan memperolehi frasa- frasa kunci, konsep dan sebagainya, dan persediaan untuk pemprosesan teks dengan cara yang betul untuk analisis yang seterusnya dengan teknik olahan data numerik seperti untuk mengenalpasti “co-occurrences” bagi konsep, nama, alamat, nama produk dan sebagainya.

Olahan teks menggunakan “Precision” dan “Recall” untuk mengukur keberkesanan pelbagai teknik pemerolehan maklumat yang membenarkan perbandingan kuantitatif dilakukan. Menurut Wolfgang (2003), terdapat dua cara untuk mengukur pemerolehan olahan teks iaitu “Precision” dan “Recall”. “Precision” digunakan untuk mengenalpasti bilangan dokumen-dokumen yang diperolehi semula daripada dokumen-dokumen yang ada yang merupakan fakta

fakta yang betul. “Recall” pula digunakan untuk mengenalpasti bilangan dokumen-dokumen yang sepatutnya diperolehi semula yang berada di dalam fakta perolehan semula. Berikut adalah merupakan formula-formula yang digunakan untuk mengukur “Precision” dan “Recall” berdasarkan Wolfgang (2003) :

$$\text{Precision} = \frac{\text{Relevant \& Retrieved}}{\text{Retrieved}}$$
$$\text{Recall} = \frac{\text{Relevant \& Retrieved}}{\text{Relevant}}$$

Rajah 2.1 di bawah ini menerangkan secara ringkas berkaitan “Relevant documents” dan “Retrieved documents” :



Rajah 2.1 : Perkaitan antara “Relevant documents” dan “Retrieved documents”

Menurut Li dan Yamanishi (2001), daripada kajian yang telah mereka lakukan, mereka telah mendapatkan 10 katakunci daripada jawapan-jawapan imej untuk setiap jenis kereta. Seterusnya pengujian dilakukan bagi mengenalpasti takat atau had yang dipersetujui oleh setiap katakunci dengan perkataan-perkataan yang wujud di dalam 10 peraturan bagi “classification” dan “association” bagi output yang dihasilkan oleh SA (“Survey Analyzer”). Proses ini dilakukan dengan menggunakan “Precision” dan “Recall”.

Di sini, “Precision” ditakrifkan sebagai nisbah bagi nombor perkataan-perkataan yang betul yang diperolehi kepada jumlah nombor bagi perkataan-perkataan yang diperolehi. Manakala “Recall” pula ditakrifkan sebagai nisbah bagi nombor perkataan-perkataan yang betul yang diperolehi kepada jumlah nombor bagi perkataan-perkataan yang akan diperolehi. Jadual 2.2 menunjukkan Keputusan Penilaian bagi “Precision” dan “Recall” yang diambil daripada sumber Li dan Yamanishi (2001).

| | Association Rule | | Classification Rule | |
|----------------|------------------|-----------|---------------------|-----------|
| | Recall | Precision | Recall | Precision |
| Car A | 0.80 | 0.90 | 0.70 | 0.80 |
| Car B | 0.90 | 0.90 | 0.90 | 0.90 |
| Car C | 0.90 | 0.90 | 0.80 | 0.80 |
| Car D | 0.50 | 1.00 | 0.40 | 1.00 |
| Car E | 0.80 | 0.80 | 0.70 | 1.00 |
| Car F | 0.60 | 0.90 | 0.60 | 1.00 |
| Average | 0.75 | 0.90 | 0.68 | 0.92 |

Jadual 2.2 : Keputusan Penilaian

2.4 Soalan Terbuka

Soalan terbuka membolehkan responden memberikan jawapan secara bebas dengan memberikan jawapan mereka sendiri kepada soalan terbuka yang diberikan. Bagi kajian ini, setiap jawapan kepada soalan terbuka yang diedarkan kepada pelajar dijadikan sebagai input kepada sistem yang dibangunkan. Soalan terbuka ini meminta agar pelajar menyatakan ciri-ciri yang mereka suka pada telefon bimbit yang mereka miliki. Terdapat empat jenis telefon bimbit sahaja yang dipilih iaitu Nokia, Motorola, Siemen dan Samsung. Bilangan responden bagi setiap jenis telefon bimbit ialah masing-masing sebanyak sepuluh orang. Sila rujuk lampiran untuk contoh soalan terbuka yang digunakan sebagai input kepada Sistem Olahan Teks Dengan Menggunakan Soalan Terbuka ini.

2.5 “Classification Rule”

Menurut Yi dan Yamanishi (2001), “classification rule” mengandungi satu rangkaian atau siri mengikut tertib atau turutan peraturan “IF-THEN-ELSE” untuk menentukan atau menetapkan jawapan terbuka kepada sasaran. Setiap peraturan ini mempunyai satu syarat atau keadaan untuk penentuan yang memerlukan kewujudan secara serentak beberapa perkataan atau kewujudan satu perkataan. Setiap peraturan juga disertakan dengan nilai kebarangkalian atau kemungkinan (frekuensi relatif) kepada penentuannya.

Li dan Yamanishi (2001), telah membuat kajian yang bertajuk “Mining Open Answers in Questionnaire Data” yang melibatkan jenis kereta sebagai sasaran dan imej bagi jenama kereta sebagai jawapan terbuka. Jadual 2.3 menunjukkan contoh bagi data soal-selidik yang menunjukkan jenis kereta dan input-input yang diperolehi.

| Car | Brand Image |
|------------|---------------------|
| Car A | For ordinary people |
| Car A | Easy to drive |
| | |
| Car B | High performance |
| Car B | Mobility |
| | |

Jadual 2.3 : Contoh data soal selidik

Jadual 2.4 pula menunjukkan histogram perkataan untuk imej kereta jenama A sebagai output bagi “Survey Analyzer” atau singkatannya, SA.