

# PENINGKATAN KINERJA CLUSTERING DOKUMEN TEKS MENGUNAKAN PEMBOBOTAN SAMPEL

**Amir Hamzah**

Jurusan Teknik Informatika, Fakultas Teknologi Industri,  
Institut Sains & Teknologi AKPRIND, Jalan Kalisahak 28, Yogyakarta 55222  
Phone: (0274)-563029; Fax: (0274)-563847  
e-mail: [miramzah@yahoo.co.id](mailto:miramzah@yahoo.co.id)

## Abstrak

Algoritma clustering berbasis pembobotan sampel (*sample weighting*) saat ini banyak diteliti. Ada beberapa model pembobotan yang pada prinsipnya bertujuan untuk merubah nilai vektor sampel dan formula similaritas vektor sampel dengan pusat klusternya. Dalam dokumen teks pembobotan dapat berupa konektivitas antar dokumen, misalnya dalam dokumen akademik yang ada koneksi referensi. Namun dalam dokumen berita koneksi referensi mungkin jarang ditemukan. Dalam makalah ini teknik pembobotan baru diajukan, yaitu menggunakan kata-kata yang muncul dalam kata kunci (*keyword*) dan judul (*title*) dari suatu dokumen teks. Eksperimen dilakukan terhadap abstrak dokumen akademik sebanyak 500 dokumen dan dokumen berita. Sebanyak 3000 dokumen Algoritma yang diuji kinerjanya adalah algoritma *K-Means clustering* dan algoritma *Fuzzy C-Means clustering*. Parameter kinerja algoritma digunakan nilai *F-measure* dari hasil clustering sebelum dilakukan pembobotan sampel dan setelah dilakukan pembobotan sampel. Hasil eksperimen menunjukkan bahwa pembobotan sampel dapat meningkatkan kinerja clustering sebesar 12,8% untuk pembobotan dengan *keyword* dan *title* dan meningkatkan kinerja clustering 9,8% untuk pembobotan dengan *title* saja.

**Kata kunci:** *clustering dokumen, pembobotan sampel, kinerja clustering, F-measure*

## 1. PENDAHULUAN

Analisis kluster merupakan alat analisis statistik multivariate dan merupakan salah satu metode penting dalam pengenalan pola tak terbimbing (*unsupervised pattern recognition*). Aplikasi analisis ini sangat luas dan terutama peranan pentingnya dalam konteks penambangan data (*data mining*). Sebagai teknik *unsupervised* metode clustering dapat dikelompokkan menjadi beberapa pendekatan, antara lain *partitional clustering*, *hierarchical clustering*, *density-based clustering*, *grid-based clustering* dan *model-based clustering* (Han and Kamber, 2000). Saat ini, dengan semakin melimpahnya data digital yang dihasilkan oleh jaringan komputer internet, analisis kluster menjadi alat analisis yang handal.

Dalam bidang Sistem Temu Kembali Informasi (*Information Retrieval System*), metode clustering juga telah diterapkan pada berbagai sisi, misalnya dalam mempartisi corpus (Grossman and Frieder, 2004), mengekstrak konsep (Karypis, 2000), atau meningkatkan kinerja clustering dengan membangun Sistem Temu Kembali berbasis konsep (Hamzah, 2009).

Salah satu kesulitan dalam clustering dokumen teks dengan model ruang vektor berbasis kata adalah bermula dari asumsi bahwa kata-kata dalam dokumen saling independen sedemikian sehingga perhitungan jarak antar dokumen yang diwakili oleh jarak antar vektor dokumen dalam ruang vektor dapat ditetapkan menggunakan berbagai formula jarak. Jika asumsi ini tidak dipenuhi maka perhitungan jarak sebenarnya menjadi kurang akurat. Meskipun pada clustering dokumen ukuran kedekatan lebih sering digunakan ukuran similaritas dari pada fungsi jarak, tetapi efek tidak terpenuhinya asumsi independensi tetap terjadi. Pada kenyataannya lebih sering antar kata dalam suatu dokumen adalah tidak independen, justru kata yang satu terkait secara makna dengan kata yang lain.

Persoalan lain selain independensi kata sebagai dimensi dalam ruang vektor adalah dalam algoritma yang ada, misalnya *K-Means clustering* atau *Fuzzy C-Means clustering* juga diasumsikan bahwa setiap sampel dalam suatu kluster dianggap memiliki bobot yang sama pentingnya terhadap prototype atau pusat kluster. Pada kenyataannya sebenarnya tidak demikian karena sesuai dengan similaritas yang berbeda antara dokumen dalam kluster dengan pusat klusternya maka mestinya ketika mengukur similaritas juga menggunakan bobot formula yang berbeda. Menyadari hal ini telah dicoba beberapa pendekatan untuk memberikan bobot yang berbeda pada fungsi similaritas dokumen atau fungsi tujuan dalam clustering, antara lain yang diajukan oleh Zhang and Zhou (2006) yang menggunakan hubungan referensi (*cited relationship*) antara dokumen sebagai pembobot algoritma. Li et.al. (2005) dan Bao et.al. (2006) memadukan mean dan modus untuk memberi bobot dalam algoritma *Fuzzy C-Means*. Kelemahan pembobotan menggunakan hubungan referensi adalah keterbatasan algoritma hanya pada penerapan clustering dokumen akademik yang memang lazim memiliki data referensi. Pada jenis dokumen lain seperti dokumen berita atau dokumen akademik tetapi hanya berupa abstrak dokumen, maka informasi referensi tidak ada, sehingga algoritma ini tidak dapat diterapkan. Oleh karena itu dalam penelitian ini diusulkan

pembobotan baru dengan menggunakan informasi dalam judul dokumen dan kata kunci dokumen dalam abstrak akademik sebagai faktor pembobot dalam algoritma.

## 2. TINJAUAN PUSTAKA

### Algoritma Clustering dokumen menggunakan pembobotan sample.

Latar belakang ide dibalik algoritma pembobotan sampel adalah asumsi bahwa sampel atau objek yang berbeda memiliki peran yang berbeda dalam proses clustering. Dengan demikian harus diupayakan agar bobot yang berbeda diberikan pada objek atau sample yang berbeda. Clustering dengan pembobotan sample ini dapat meningkatkan kinerja clustering (Nock and Frank, 2006).

### Pembobotan Sampel pada Algoritma K-Mean.

Jika tidak melibatkan pembobotan sampel, algoritma klasik K-Means iterasinya berakhir ketika fungsi tujuan telah konvergen ke suatu nilai yang tetap. Fungsi tujuan dalam algoritma clustering dokumen dapat dituliskan sebagai :

$$J = \sum_{i=1}^K \sum_{j=1}^{m_i} sim(\bar{d}_j, \bar{c}_i) \quad (1)$$

$$\bar{c}_i = \frac{1}{m_i} \sum_{j=1}^{m_i} \bar{d}_j \quad (2)$$

Dengan J akan konvergen pada suatu nilai apabila iterasi dalam clustering tidak lagi berubah komposisi sampelnya pada setiap kluster. Nilai J dapat digunakan untuk mengukur kinerja clustering. Jika menggunakan rumus jarak, maka semakin kecil nilai J semakin baik kinerja clusteringnya, karena berarti sampel-sampel dalam kluster mendekati pusat klusternya secara optimal. Jika menggunakan rumus similaritas, semakin besar nilai J semakin baik hasil clustering. Adapun K adalah banyaknya kluster, sedangkan  $m_i$  adalah banyaknya sampel dalam kluster ke-i. Adapun  $\bar{d}_j$ , adalah vektor objek ke-j dalam kluster ke-i, sedangkan  $\bar{c}_i$  adalah vektor pusat kluster ke-i yang dapat dihitung menggunakan rumus (2).  $sim(\bar{d}_j, \bar{c}_i)$  adalah tingkat similaritas antara vektor pusat kluster ke-i dengan sampel dokumen ke-j. Keduanya diwakili oleh model ruang vektor setelah ekstraksi fitur dan komputasi pembobotan fitur dilakukan. Dengan demikian setiap dokumen akan diwakili oleh vektor dokumen yang berisi elemen vektor  $w_{ij}$  yang bernilai real dan mewakili tingkat kepentingan fitur (kata) ke i dalam dokumen ke-j. Pembobotan yang biasa digunakan untuk pembobotan fitur dalam clustering dokumen adalah TF-IDF ternormalkan. Formula TF-IDF yang digunakan adalah sebagai rumus (3) berikut (Chisholm dan Kolda, 1999) :

$$w_{ij} = \frac{(\ln(f_{ij} + 1) \cdot \log\left(\frac{N}{n_i}\right))}{\sqrt{\sum_{i=1}^I \left( (\ln(f_{ij} + 1) \cdot \log\left(\frac{N}{n_i}\right)) \right)^2}} \quad (3)$$

Pengukuran fungsi similaritas digunakan similaritas cosine (Salton and McGill, 1983). Fungsi cosine merupakan fungsi terbaik dari sisi kemampuan mengukur kesamaan objek dan efisiensi komputasi untuk clustering dokumen dibandingkan dengan jarak *eulidean*, fungsi *jaccard*, fungsi *dice* atau korelasi *pearson* (Hamzah, 2008).

Dalam algoritma clustering terbobot fungsi tujuannya adalah seperti pada persamaan (4) berikut :

$$J' = \sum_{i=1}^K \sum_{j=1}^{m_i} (w_j \bullet sim(\bar{d}_j, \bar{c}_i)) \quad (4)$$

dimana  $w_j$  adalah bobot untuk sampel ke-j dengan kendala yang harus dipenuhi adalah  $\sum_{j=1}^{m_i} w_j = 1$  dan  $\bar{c}_i$

adalah prototipe (pusat kluster) dari kluster ke-i setelah proses clustering dengan sampel yang diberi bobot, yang dapat dihitung menggunakan rumus :

$$\bar{c}_i = \sum_{j=1}^{m_i} w_j \bullet \bar{d}_j \quad (5)$$

Terlihat dari rumus (5) bahwa bobot  $w_j$  memerankan peranan penting dalam proses penyesuaian prototipe dari suatu kluster. Jika nilai ini diberi bobot yang sama untuk setiap vektor sampel  $\bar{d}_j$ , misalnya nilainya adalah  $1/m_i$ , dengan nilai  $m_i$  adalah cacah sampel dalam kluster ke- $i$ , maka prototype atau pusat kluster ke- $i$ , yaitu  $c_i$  akan berubah menjadi pusat kluster biasa tanpa pembobotan seperti pada persamaan (2). Sebagai akibatnya juga rumus persamaan (4) akan menjadi persamaan (1). Jika bobot  $w_j$  diambil nilai sama pada setiap nilai vector sampel  $\bar{d}_j$  maka algoritma akan berubah menjadi algoritma K-Means clustering biasa tanpa pembobotan.

### Pembobotan Sampel pada Algoritma Fuzzy C-Mean

Formula iterative dalam Fuzzy C-Means untuk clustering dokumen adalah sebagai rumus persamaan (6) ini:

$$u_{ij} = \frac{\text{sim}^2(\bar{d}_j, \bar{c}_i)}{\sum_{k=1}^c \text{sim}^2(\bar{d}_j, \bar{c}_k)} \quad (6)$$

$$c_i = \frac{\sum_{j=1}^n u_{ij}^2 \bar{d}_j}{\sum_{j=1}^n u_{ij}^2} \quad (7)$$

dengan nilai  $u_{ij}$  adalah bobot keanggotaan sampel ke- $j$  pada cluster ke- $i$ , nilai  $c$  adalah cacah kluster. Pada setiap iterasi nilai pusatkluster di-update menggunakan rumus persamaan (7) setelah sebelumnya nilai keanggotaan sampel ke- $j$  dalam kluster ke- $i$  diupdate dengan persamaan (6). Iterasi dilakukan sampai fungsi tujuan bernilai relative tetap. Adaoun fungsi tujuan dalam Fuzzy C-Means adalah seperti persamaan (8) berikut :

$$J = \sum_{i=1}^c \sum_{j=1}^n (u_{ij}^2 \cdot \text{Sim}^2(\bar{d}_j, \bar{c}_i)) \quad (8)$$

Dalam konteks clustering Fuzzy C-Means dengan pembobotan sampel maka fungsi criteria pada persamaan (8) berubah menjadi sebagai berikut (persamaan 89) :

$$J' = \sum_{i=1}^c \sum_{j=1}^n w_j \cdot (u_{ij}^2 \cdot \text{Sim}^2(\bar{d}_j, \bar{c}_i)) \quad (9)$$

Pusat cluster pada persamaan (7) jika clustering dilakukan dengan pembobotan sampel akan berubah menjadi seperti persamaan (10) berikut :

$$c_i = \frac{\sum_{j=1}^n w_j \cdot u_{ij}^2 \bar{d}_j}{\sum_{j=1}^n w_j \cdot u_{ij}^2} \quad (10)$$

### Metode Pembobotan Sampel

Pembobotan sampel untuk *K-Means Clustering* (persamaan (4) dan (5)) dan *Fuzzy C-Means Clustering* (persamaan (9) dan (10)) dilakukan untuk dengan memberi bobot pada dokumen dengan mengambil nilai bobot TF-IDF kata-kata yang muncul dalam title, atau dalam kata kunci. Bobot untuk title dalam suatu dokumen, dicatat sebagai  $w_j^{\text{TITLE}}$  adalah rasio antara total bobot TF-IDF kata yang berasal dari title dengan total bobot sejumlah  $n$  kata terpenting dalam koleksi, dengan  $n$  ditetapkan berdasarkan eksperimen. Sedangkan nilai bobot untuk kata kunci, ditulis sebagai  $w_j^{\text{KEYWORD}}$  adalah rasio antara total bobot TF-IDF kata yang berasal dari kata kunci dengan total bobot sejumlah  $n$  kata terpenting dalam koleksi. Formula dari kedua bobot tersebut dapat ditulis dalam persamaan (11) dan (12) berikut :

$$w_j^{\text{TITLE}} = \frac{\sum_{k=1}^T w_{jk}}{\sum_{k=1}^n w_k} \quad (11)$$

$$w^{\text{KEYWORD}}_j = \frac{\sum_{k=1}^T w_{jk}}{\sum_{k=1}^n w_k} \quad (12)$$

### 3. METODE PENELITIAN

Dalam penelitian ini digunakan bahan objek atau sampel data berupa koleksi dokumen teks, yang berupa dokumen teks berita dan dokumen teks abstrak dari suatu makalah ilmiah. Secara statistik koleksi dokumen memiliki karakteristik sebagai berikut :

**Tabel 1.** Koleksi dokumen teks yang dijadikan percobaan clustering

Koleksi	Cacah dok	Frek min kata	Frek mak kata	Rata2 Frek kata
News3000	3000	358	2254	998
Abs500	500	198	324	198

Adapun format dokumen telah dipersiapkan sedemikian sehingga mendukung proses clustering dengan memformat menjadi tiga bagian untuk dokumen akademik abstrak, yaitu <TITLE></TITLE>, <BODY></BODY> dan <KEYWORD></KEYWORD>. Untuk dokumen abstrak formatnya adalah sebagai Gambar 1 berikut:

```

<DOC>
<DOCNO>abs-006</DOCNO>
<TITLE> perangkat lunak berbasis SMS</TITLE>
<BODY> pengembangan perangkat lunak sistem kendali
dan pengawasan menggunakan relay on off berbasis sms
dan database untuk data historis Pada penelitian ini
akan dikembangkan suatu perangkat lunak bantu
pengontrol dan monitor suatu keamanan ruangan
berbasis Short Message Service (SMS) dan...
</BODY>
<KEYWORD>SMS,perangkat lunak,database</KEYWORD>
</DOC>
    
```

**Gambar 1** Format koleksi dokumen abstrak

Untuk dokumen berita formatnya tidak banyak berbeda dengan format dokumen akademik abstrak, bedanya adalah dokumen berita tidak memiliki keyword , sehingga tidak ada tag <KEYWORD></KEYWORD>.

Pembuatan program dilakukan dengan menggunakan komputer PC Intel Pentium IV 2.8GHz, RAM 1GB, Hard Disk 80 GB, dan sistem operasi Windows XP Professional. Bahasa pemrograman yang dipergunakan adalah java jdk1.6.4, dan Matlab versi 7.0.4.

Algoritma K-Means clustering yang dikodekan dan diuji adalah sebagai berikut :

- [1] Ambil K objek sebagai seed dari K pusat kluster
- [2] Untuk semua objek: cari kluster dengan jarak terdekat, dan tetapkan objek masuk dalam kluster tersebut.
- [3] Hitung ulang pusat kluster dengan rata-rata objek dalam kluster tersebut
- [4] Hitung fungsi kriteria dan lakukan evaluasi. Jika fungsi kriteria berubah cukup kecil algoritma berhenti.

Algoritma tersebut sebenarnya sama untuk K-Means tanpa pembobotan, bedanya adalah dalam hal perhitungan pusat kluster dan fungsi tujuan. Langkah 3 dilakukan dengan menggunakan formulas rata-rata sampel terbobot seperti persamaan (5), sedangkan langkah 4 dilakukan dengan menggunakan fungsi tujuan terbobot seperti persamaan (4).

Algoritma Fuzzy C-Means clustering yang dikodekan dan diuji adalah sebagai berikut :

- [1] Tentukan secara random matrik U inisial yang beranggotakan  $u_{ij}$ , sebagai bilangan pecah real sehingga memenuhi kondisi :  $\sum_{j=1}^C u_{ij} = 1$
- [2] Update pusat kluster menggunakan persamaan (10)
- [3] Update matrik U menggunakan persamaan (6)
- [4] Hitung fungsi tujuan menggunakan persamaan (9)
- [5] Jika perubahan nilai fungsi tujuan sudah lebih kecil dari suatu bilangan kecil tertentu, algoritma berhenti, jika tidak ulangi langkah [2].

### Evaluasi kinerja clustering

Untuk evaluasi hasil clustering, atau dikenal sebagai validitas clustering dapat dilakukan secara eksternal atau internal. Validitas eksternal adalah menguji sejauh mana algoritma mampu merekonstruksi objek sampel yang sebelumnya sudah dikategorikan dengan label atau class tertentu. Langkah untuk validitas *clustering* eksternal adalah dengan menyusun matriks konfusi (*confusion matrix*), yaitu matriks yang disusun berdasarkan berapa banyak objek yang diklasifikasikan dengan benar oleh proses *clustering*.

### Entropy

Jika  $p_{ij}$  adalah peluang anggota kluster  $j$  adalah milik klas  $i$ , maka *entropy* tiap kluster ditentukan sebagai :

$$E_j = - \sum_i p_{ij} \log(p_{ij}) \quad (13)$$

dan total *entropy* untuk hasil suatu *clustering* pada *predifined-class objects* adalah :

$$E = \sum_j \frac{n_j * E_j}{n} \quad (14)$$

di mana  $n_j$  = cacah objek dalam kluster  $j$ ,  $m$  = cacah kluster dan  $n$  = total objek.

### F-measure

*F-measure* merupakan kombinasi ide *precision* dan *recall* dari *information retrieval* (Rijsbergen, 1979). Untuk kluster  $j$  dan klas  $i$  didefinisikan :

$$R = Recall(i, j) = n_{ij}/n_i$$

$$P = Precision(i, j) = n_{ij}/n_j$$

*F-measure* untuk klas  $i$  adalah :

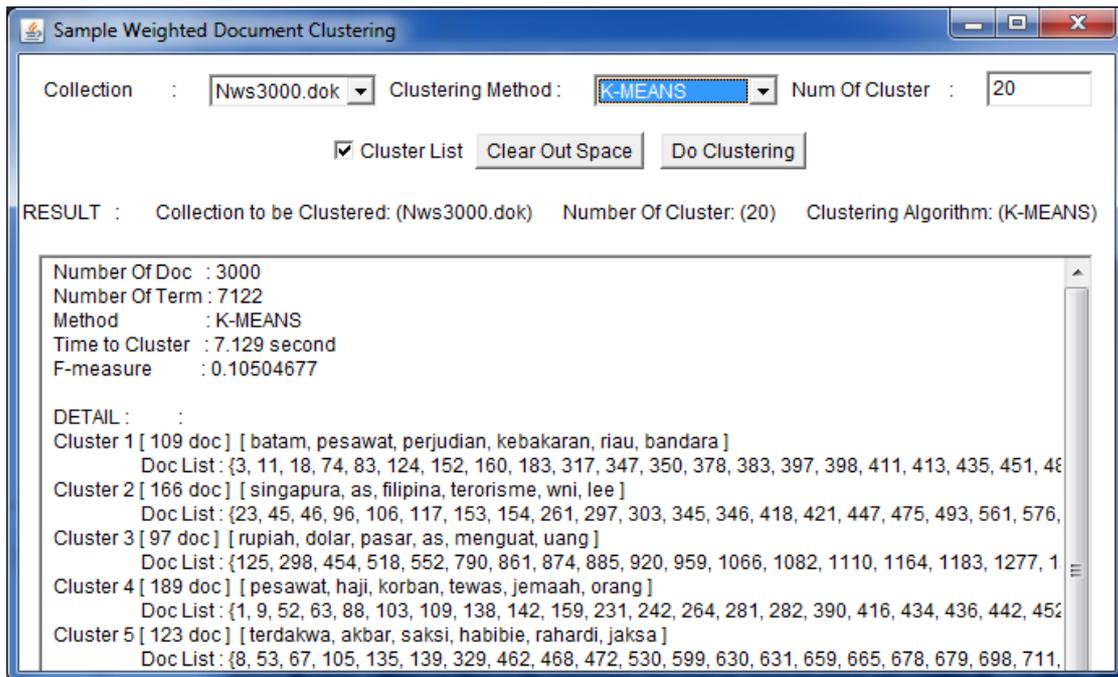
$$F(i) = 2PR/(P+R) \quad (15)$$

di mana  $F(i)$  diambil nilai terbesar dari setiap kluster untuk klas  $i$ . *F-measure* keseluruhan kluster hasil *clustering* adalah :

$$F = \frac{\sum_i n_i * F(i)}{n} \quad (16)$$

## 4. HASIL DAN PEMBAHASAN

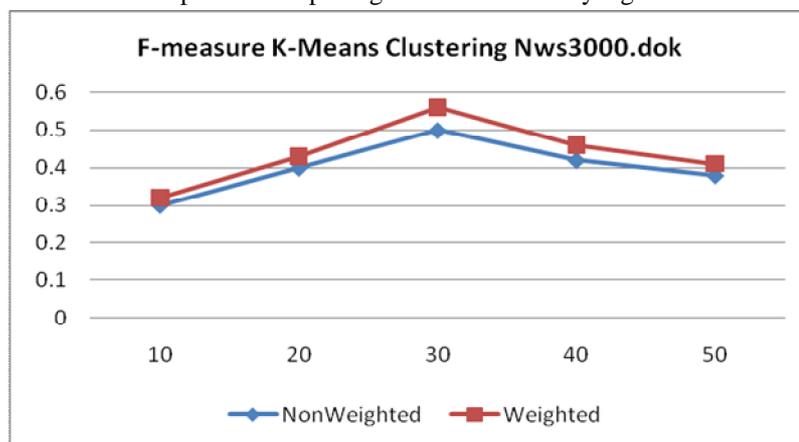
Hasil perancangan program *Sample Weighted Document Clustering* disajikan seperti contoh output program pada Gambar 2. Pada hasil clustering tersebut terlihat metode clustering dengan K-MEANS pada koleksi dokumen berita Nws3000.dok dengan jumlah kluster 20 kluster. Proses kluster diselesaikan dalam waktu 7.129 detik, dengan nilai *F-measure* 0.10504677.



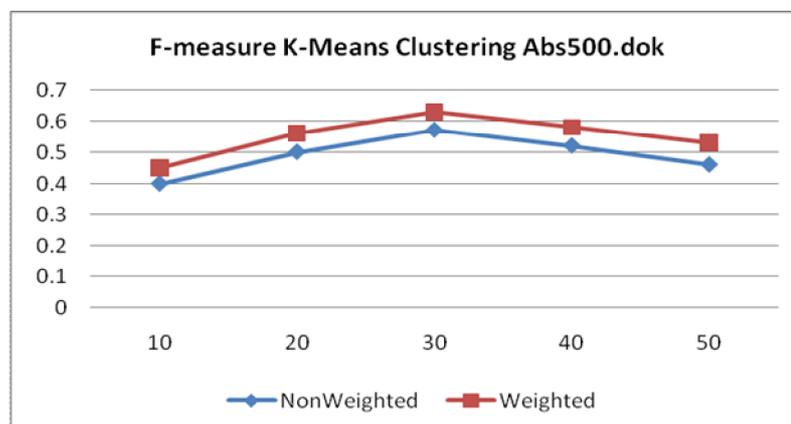
Gambar 2. Antar Muka Program Sample Weighted Clustering

### Perbandingan clustering pada algoritma K-Means

Berikut ini disajikan perbandingan kinerja clustering pada algoritma K-Means clustering, pada koleksi dokumen berita berdasarkan nilai F-measure pada beberapa tingkatan nilai kluster yang dicobakan.



Gambar 3. Nilai F-measure Hasil K-Mean Clustering pada Dokumen Berita

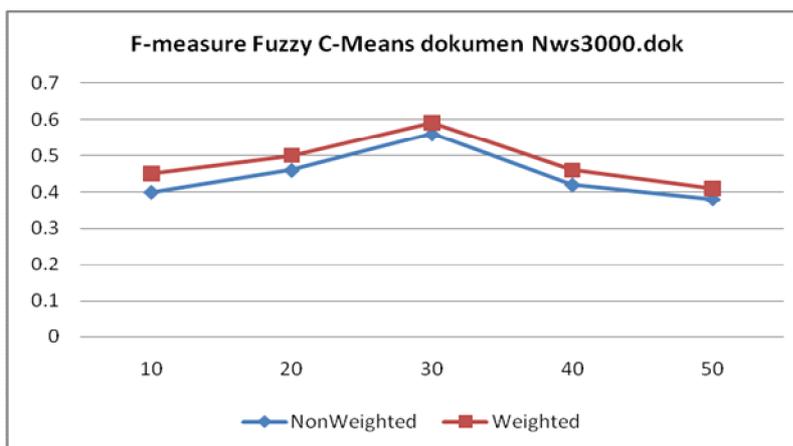


Gambar 4. Nilai F-measure Hasil K-Means Clustering pada okumen Abstrak

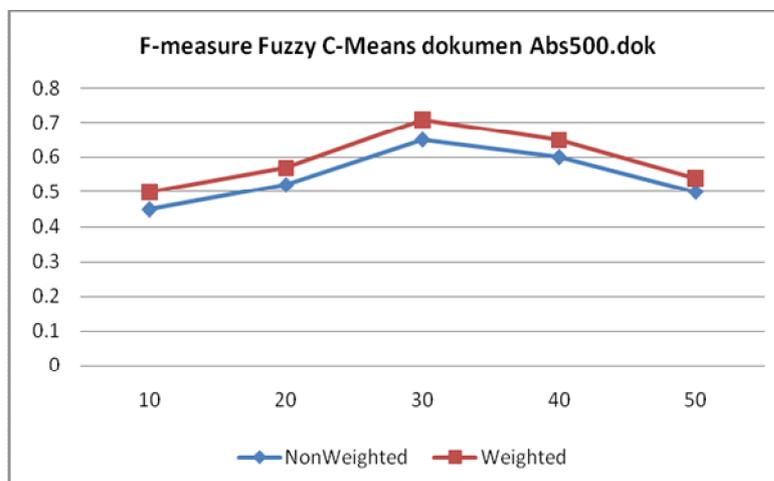
Dari gambar 3 dan gambar 4 terlihat bahwa pemberian bobot pada dokumen berita dengan melibatkan faktor title sebagai bobot telah memperbaiki kinerja rata-rata sebesar 9,6% pada perbaikan nilai validitas clustering F-measurenya pada koleksi dokumen berita. Sedangkan pada dokumen abstrak, dengan melibatkan faktor title dan kata kunci sebagai bobot telah menaikkan kinerja clustering rata-rata 12,5% pada nilai F-measurenya.

### Perbandingan clustering pada algoritma Fuzzy C-Means

Berikut ini disajikan perbandingan kinerja clustering pada algoritma Fuzzy C-Means clustering, pada koleksi dokumen berita pada Gambar 5 dan dokumen abstrak pada Gambar 6.



Gambar 5. Nilai F-measure clustering Fuzzy C-Means Koleksi Dokumen Berita



Gambar 6. Nilai F-measure clustering Fuzzy C-Means Koleksi Dokumen Abstrak

Dari gambar 5 dan gambar 6 juga terlihat bahwa dalam algoritma Fuzzy C-Means juga memberikan hasil dengan pola yang sama dengan algoritma K-Means clustering dengan sedikit peningkatan kinerja dalam rata-rata F-measure, yaitu sebesar 10% untuk dokumen berita dan 12,9% untuk dokumen abstrak.

## 5. KESIMPULAN

Dari pembahasan dapat diambil beberapa kesimpulan dari penelitian ini, antara lain bahwa pemberian bobot dengan melibatkan judul dokumen dan kata kunci mampu meningkatkan kinerja clustering sehingga hasil clustering lebih baik dibandingkan dengan clustering tanpa pembobotan sampel. Hal ini terjadi baik menggunakan algoritma *K-Means clustering* maupun menggunakan algoritma *Fuzzy C-Means clustering*. Perbaikan kinerja dengan melibatkan bobot dengan menggunakan dua bobot yaitu judul dan kata kunci dari dokumen mampu meningkatkan kinerja yang lebih tinggi (sekitar 12,8%) dibandingkan dengan menggunakan

judul saja, sedangkan pembobotan hanya melibatkan judul dokumen mampu meningkatkan kinerja sekitar rata-rata 9,8% lebih tinggi dari clustering dokumen tanpa pembobotan.

#### DAFTAR PUSTAKA

- Bao, Z., Han, B., and Wu, S., 2006, *A General Weighted Fuzzy Clustering Algorithm*, Lecture Notes in Computer Science, Volume 4142/2006, 102-109, DOI:10.1007/11867661\_10.
- Chisholm, E. and T. G. Kolda, 1999, *New Term Weighting Formula for the Vector Space Method in Information Retrieval*, **Research Report**, Computer Science and Mathematics Division, Oak Ridge National Library, Oak Ridge, TN 3781-6367, March 1999.
- Hamzah, A, A. Susanto, F. Soesianto, J.E. Istiyanto, 2008, *Studi Kinerja Fungsi-Fungsi Jarak Dan Similaritas Dalam Clustering Dokumen Teks Berbahasa Indonesia*, Seminar Nasional Informatika, Prosiding Seminar Nasional SEMNASIF2008, Universitas Pembangunan Nasional "Veteran", Yogyakarta 24 Mei 2008
- Hamzah, A, 2009, *Penerapan Clustering Dokumen untuk Meningkatkan Efektifitas Sistem Temu Kembali Informasi Dokumen Berbahasa Indonesia*, **Disertasi**, Fakultas Teknik, Universitas Gadjah Mada, Yogyakarta.
- Han, J., and Kamber, M., 2000, *Data Mining: Concept and Techniques*, Morgan Kaufman.
- Grossman, D. A. and O. Frieder, 2004, *Information Retrieval Algorithms and Heuristics*, Springer, 2<sup>nd</sup> edition, 2004.
- Li, Jie, Gao, X., and Jiao, L., 2005, *A Novel typical-Sample-Weighted Clustering Algorithm for Large Data Sets*, LNAI3801, 696-703
- Karypis, G. and Han Eui-Hong, 2000, *Concept Indexing A Fast Dimensionality Reduction Algorithm with Applications to Document Retrieval and Categorization*, Technical Report TR-00-0016, University of Minnesota. [www.cs.umn.edu/karypis](http://www.cs.umn.edu/karypis)
- Nock, R., and Nielsen, F., 2004, *An Abstract Weighting Framework for Clustering Algorithms*, in: Proceedings of the Fourth International SIAM Conference on Data Mining, 200-209.
- Rijsbergen, C. J., 1979, *Information Retrieval*, Information Retrieval Group, University of Glasgow, UK
- Salton, G. and McGill, M.C., 1983, *Introduction to Modern Information Retrieval*, McGraw-Hill Book, Co., New York.
- Zhang, C., Su, Z., and Zhou, D., 2006, *Document Clustering Using Sample Weighting*, Nanjing University of Science & Technology (NO.JGQN0701).