

KOMPARASI ALGORITMA UNTUK KLASIFIKASI HEREGISTRASI CALON MAHASISWA

Dadang Aribowo¹⁾, Aris Ekyanto Heru Setiadi²⁾, Ivandari³⁾

¹⁾ Progran Studi Teknik Informatika, STMIK Widya Pratama Pekalongan

¹⁾ email: dadang.stmik.wp@gmail.com

²⁾ email: aris@stmik-wp.ac.id

³⁾ email: ivandarialkaromi@gmail.com

Abstract

Students are the most valuable asset in a private university (PTS). Because most of PTS's revenues and operating costs are obtained from students. The number of students who do registration clearly will be a breath of fresh air for the institution. In the last 5 years, around 20% of STMIK Widya Pratama students did not register. Early knowledge of prospective students who might not register will be a reference for the institution to take action to maintain students. The recording of student data that is neatly arranged can be used by management to analyze the characteristics and causes of students not registering. Data mining can process past data into new information or knowledge. In data mining, there is one major function, namely the classification that processes training data to calculate new data / data testing. Methods or algorithms that can be used in the classification process are numerous with various characteristics of each. Some of the best classification algorithms include naive bayes, knn, and C4.5. The results showed that the three algorithms, namely, naive bayes and the C45 decision tree can be used to classify prospective student registrations. The accuracy of the C45 decision tree algorithm is the best, 80.72% followed by the algorithm with an accuracy rate of 80.46%. While the accuracy of naive bayes is the lowest with 74.49%.

Keywords: KNN, Naive Bayes, Decission Tree C45

1. PENDAHULUAN

Mahasiswa merupakan aset dalam sebuah perguruan tinggi. Kemajuan perguruan tinggi dapat dilihat dari jumlah mahasiswa yang melakukan pendaftaran pada setiap pembukaan tahun ajaran baru. Perguruan tinggi favorit tentunya akan lebih diminati calon mahasiswa dibandingkan dengan perguruan tinggi yang lain. Artinya banyaknya calon mahasiswa baru menjadi salah satu indikasi kemajuan perguruan tinggi. Penerimaan mahasiswa baru di STMIK Widya Pratama dilakukan setiap tahun untuk keempat program studi yang ditawarkan. Tahapan dalam penerimaan Mahasiswa baru adalah sebagai berikut: isi formulir pendaftaran, tes ujian masuk, pengumuman hasil tes dan heregistrasi. Dalam 5 tahun trakhir selalu ada selisih yang cukup besar dari jumlah pendaftar dengan jumlah mahasiswa yang melakukan heregistrasi. Tabel 1 berikut menunjukkan data jumlah pendaftar pada 5 tahun terakhir.

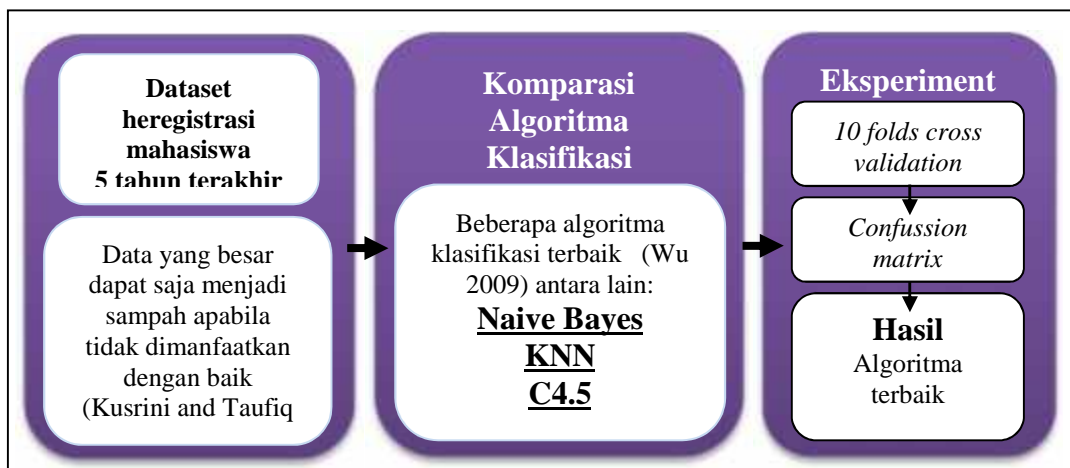
Tabel 1. Data Penerimaan Mahasiswa Baru 5 tahun terakhir

Tahun	Pendaftaran	Registrasi	Selisih	Prosentase	
				Registrasi	Tidak Registrasi
2013	705	513	192	73%	27%
2014	608	474	134	78%	22%
2015	667	500	167	75%	25%
2016	654	427	227	77%	23%
2017	506	428	78	84,6%	15,4%

Dari data tabel 1 terlihat adanya prosentase yang cukup besar yang terdapat pada calon mahasiswa yang tidak melakukan heregistrasi. Penurunan jumlah mahasiswa jelas akan berpengaruh dalam pendapatan lembaga, mengingat STMIK Widya Pratama merupakan kampus swasta. Adanya selisih jumlah yang cukup banyak antara jumlah pendaftar dan jumlah mahasiswa yang melakukan heregistrasi, membuat pihak yayasan berpikir keras agar jumlah mahasiswa STMIK Widya Pratama dapat terus dipertahankan dan ditingkatkan. Jika kemungkinan pengunduran diri mahasiswa dapat diketahui sejak dini, maka pihak manajemen dapat melakukan tindakan-tindakan preventif untuk mempertahankan calon mahasiswa tersebut (Kusrini and Taufiq 2009).

Klasifikasi adalah salah satu peran utama data mining (Witten, Frank, and Hall 2011). Proses klasifikasi dapat dilakukan dengan berbagai macam cara (Larose 2005). Berbagai macam algoritma juga dapat diaplikasikan dalam proses klasifikasi ini (Han and Kamber 2006). Berbagai masalah dalam dunia nyata banyak terselesaikan dengan menggunakan teknik klasifikasi data mining (Ashari, Paryudi, and Tjoa 2013). Penelitian klasifikasi lain juga banyak dilakukan oleh peneliti dengan menggunakan berbagai macam dataset serta algoritma yang berbeda (Ragab et al. 2014) (Patel, Vala, and Pandya 2014) (Amancio et al. 2013). Beberapa algoritma klasifikasi terbaik menurut Xindong Wu antara lain Naive Bayes, C4.5 serta K-Nearest Neighbour (Wu et al. 2007).

Klasifikasi heregistrasi mahasiswa menjadi bahan penelitian yang cukup menarik. Pada tahun 2014 Alkaromi (Alkaromi 2014) melakukan klasifikasi heregistrasi mahasiswa dengan menggunakan *K-NN*. Data yang digunakan adalah data PMB dengan *record* sebanyak 2389 dan atribut sebanyak 44. Kesemua atribut tersebut adalah atribut dasar yang didapatkan dari panitia PMB tanpa dikurangi atau dipilih sebelumnya. Untuk mengurangi jumlah atribut dalam penelitian tersebut dilakukan seleksi fitur dengan menggunakan *information gain*. Hasil dari penelitian tersebut akurasi *K-NN* yang telah dilakukan seleksi fitur menggunakan *information gain* naik hingga 83,93%. Sebelumnya algoritma *K-NN* tanpa menggunakan seleksi fitur hanya memperoleh tingkat akurasi 78,15%. Dalam penelitian tersebut digunakan alat pengukuran *confussion matrix* dan validasi yaitu *10folds cross validation*. Penelitian ini membandingkan algoritma KNN, Naive Bayes serta Decision Tree C45 untuk klasifikasi heregistrasi calon mahasiswa di STMIK Widya Pratama. Gambar 1 merupakan kerangka pemikiran dalam penelitian ini.



Gambar 1 Kerangka pemikiran penelitian

2. KAJIAN PUSTAKA

Beberapa penelitian terkait komparasi algoritma telah dilakukan dan mendapatkan hasil sebagaimana berikut:

1. Penelitian yang dilakukan Alkaromi (Alkaromi 2014) ini melakukan klasifikasi heregistrasi calon mahasiswa dengan menggunakan dataset sebanyak 2389 *record*.

Dataset yang dipakai merupakan data PMB dari tahun 2011 sampai dengan tahun 2013. Atribut yang ada dalam dataset tersebut sebanyak 44 dengan didalamnya satu atribut label dan satu atribut id. Kesemua atribut silakukan seleksi fitur dengan menggunakan *information gain* kemudian dilakukan klasifikasi dengan memanfaatkan algoritma *K-NN*. Dalam penelitian tersebut algoritma *K-NN* memperoleh tingkat akurasi sebesar 78,15% dengan menggunakan kesemua atributnya. Sedangkan setelah dilakukan seleksi fitur menggunakan *information gain* akurasi dari algoritma *K-NN* naik hingga mencapai 83,93%. Penelitian ini menggunakan *confussion matrix* untuk evaluasi serta *10folds cross validation* untuk proses validasinya.

2. Penelitian yang dilakukan Devi Sugianti (Sugianti 2012) di STMIK Widya Pratama Pekalongan menghasilkan akurasi sebesar 78% menggunakan algoritma *Bayesian Classification* untuk memprediksi heregistrasi mahasiswa. Penelitian tersebut menggunakan data PMB STMIK tahun 2011 dengan 913 *record*. Atribut yang digunakan adalah: kota asal, program studi, status daftar dan gelombang. Dalam penelitian tersebut menggunakan dua *class* dalam label atau atribut tujuan, yaitu “Registrasi” dan “Tidak Registrasi”. Perhitungan manual menggunakan rumus algoritma Bayes dari data sampel dan data baru yang belum diketahui dijelaskan secara terperinci. Namun tidak dilakukan evaluasi terhadap algoritma seperti halnya dengan menggunakan *confussion matrix*. Dalam penelitian ini juga tidak dilakukan perhitungan validasi seperti halnya menggunakan *10folds cross validation*. Sehingga tidak dapat dilihat perhitungan performa dari algoritma Bayes untuk dataset tersebut.
3. Dalam penelitian klasifikasi yang dilakukan oleh Tacbir Hendro Pudjianto, Faiza Reinaldi dan Age Teogunadi pada tahun 2011 ini digunakan algoritma ID3 (Pudjianto, Renaldi, and Teogunadi 2011). Data yang digunakan adalah data Penerimaan Mahasiswa Baru UNJANI pada 5 tahun terakhir. Atribut yang digunakan dalam penelitian tersebut adalah: Kode Progdi, Jenis Kelamin, Agama, Gol.Darah, Pekerjaan, Penghasilan, Jurusan SLTA, Asal Sekolah, Kode Progdi1, Kode Progdi2, Kode Progdi3, Gelombang, Registrasi. Atribut terakhir yaitu registrasi dijadikan sebagai atribut tujuan atau label dengan 2 kemungkinan yaitu registrasi atau tidak registrasi. Dari hasil penelitian yang telah dilakukan didapatkan tingkat akurasi dari algoritma ID3 untuk dataset tersebut diatas adalah 61,89%. Hasil dari penelitian tersebut berikutnya dijadikan sebagai acuan untuk membangun sebuah sistem pendukung keputusan yang nantinya dapat digunakan untuk pihak kampus. Sistem tersebut nantinya juga akan digunakan sebagai alat ukur ketercapaian target penerimaan mahasiswa baru yang ditetapkan oleh pihak manajemen.
4. Dalam penelitian yang dilakukan Kusrini dkk (Kusrini et al. 2009). menerapkan penggunaan algoritma *K-Nearest Neighbour* untuk kemungkinan heregistrasi mahasiswa STMIK AMIKOM Yogyakarta. Sebelumnya Kusrini juga telah melakukan penelitian yang sama dengan menggunakan algoritma C4.5. Penelitian tersebut dilakukan karena dari data sebelumnya diketahui hanya sekitar 75% dari calon mahasiswa pendaftar yang melakukan heregistrasi. Data yang digunakan dalam penelitian tersebut diambil dari panitia PMB STMIK AMIKOM Yogyakarta. Atribut yang digunakan dalam penelitian ini antara lain: NEM, Jenis Kelamin, Jurusan, Gelombang, Pilihan1, Catatan, Nilai dan Jurusan Lulus. Hasil dari penelitian ini menyatakan bahwa algoritma *K-Nearest Neighbour* tidak lebih baik dari algoritma C4.5 untuk klasifikasi heregistrasi calon mahasiswa STMIK AMIKOM Yogyakarta. Dalam proses klasifikasi menggunakan *K-Nearest Neighbour* diperlukan sebuah proses pembobotan untuk setiap atribut yang ada serta dilakukan perhitungan kedekatan seluruh data *training* dengan data *testing*. Proses tersebut jelas akan memerlukan banyak waktu untuk perhitungan terlebih lagi jika atribut yang digunakan semakin banyak serta dataset semakin besar.

3. METODE PENELITIAN

Metode penelitian yang akan digunakan dalam penelitian ini adalah eksperimental. Tahapan penelitian antara lain konsep perencanaan dengan pengumpulan data, sampai dengan validasi dan evaluasi algoritma yang akan dibahas dalam sub bab berikut:

3.1 Pengumpulan Data

Tahapan pengumpulan data dalam penelitian ini dilakukan dengan mengambil data dari panitia penerimaan mahasiswa baru selama 5 tahun terakhir. Data tersebut dikonfersi dan dijadikan menjadi sebuah dataset baru. Data mentah berisikan 44 atribut data dan hanya akan digunakan sejumlah 18 atribut untuk proses klasifikasi berikutnya. Pengurangan atribut ini dilakukan karena beberapa atribut tidak terkait dan memiliki tingkat varian yang sangat tinggi. Atribut yang tidak digunakan seperti halnya nama, alamat, nomor ktp, nomor telepon serta beberapa atribut lain. Atribut data yang akan digunakan dalam penelitian ini antara lain: kode konsentrasi, biaya kuliah, gelombang grade, shift kelas, sesi, gelombang daftar, status daftar, kode daftar, pendidikan ortu, tahun lulus, status sipil, jenis kelamin, status pekerjaan, progdi, kelas, jenjang, kota asal, serta satu atribut label yaitu status registrasi Desain Eksperimen dan Pengujian

3.2 Desain Eksperimen Algoritma

Tahapan selanjutnya setelah pengumpulan data adalah desain eksperimental algoritma dilanjutkan dengan pengujian algoritma. Dalam tahapan ini nantinya akan dibandingkan ketiga algoritma untuk klasifikasi heregistrasi mahasiswa.

3.3 Tahap Eksperimen

Tahap eksperimen dilakukan dengan menggunakan *tools software* Rapid Miner. Dalam tahapan ini nantinya akan dilakukan perhitungan terhadap dataset yang telah terkumpul sebelumnya. Proses ini memungkinkan perulangan dan akan memilih algoritma dengan tingkat akurasi terbaik. Dalam tahap eksperimen ini dilakukan pula validasi dan evaluasi terhadap ketiga algoritma klasifikasi.

3.4 Validasi

Dalam proses validasi penelitian ini akan digunakan *10 folds cross validation*. Proses ini banyak digunakan oleh peneliti karena sudah terbukti baik dan menghasilkan tingkat akurasi yang stabil. Secara tori *10 folds cross validation* sudah dijelaskan secara lebih terinci dalam bab sebelumnya (Witten, Frank, and Hall 2011).

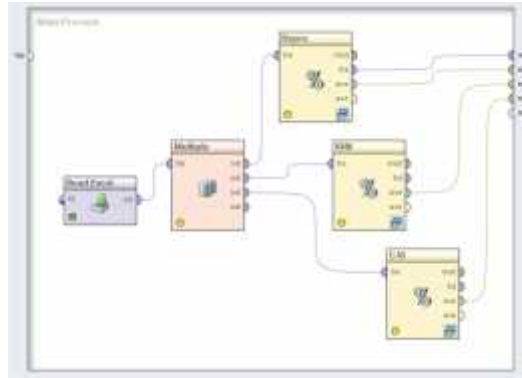
3.5 Pengukuran akurasi algoritma

Pengukuran dari suatu algoritma merupakan suatu pembuktian yang banyak dilakukan peneliti (Amancio et al. 2013). Dalam prosesnya banyak cara dapat digunakan untuk mengetahui performa suatu algoritma. Salah satu yang paling banyak digunakan adalah dengan menggunakan *confussion matrix* untuk menghitung akurasi algoritma. Perhitungan akurasi adalah presentase dari jumlah data *testing* dengan klasifikasi yang sesuai dengan aslinya dibagi keseluruhan data. Cara lain adalah dengan menghitung *Error rate*. *Error rate* adalah kebalikan dari tingkat akurasi, yaitu presentase kesalahan klasifikasi dibagi dengan keseluruhan dataset.

4. HASIL PENELITIAN

4.1 Hasil Eksperimen

Dalam tahapan eksperimen digunakan aplikasi rapid miner untuk mengolah data. Aplikasi ini dapat melakukan komparasi antara beberapa algoritma terpilih untuk melakukan perhitungan pada satu dataset yang sama. Dalam tahapan ini akan dilakukan validasi menggunakan *10 folds cross validation* serta pengukuran akurasi algoritma menggunakan *confussion matrix*. Gambar 2 merupakan desain eksperimen penelitian yang dilakukan.



Gambar 2. Tampilan desain penelitian

Dalam gambar 2 tersebut terdapat beberapa proses yang juga telah ditandai menggunakan angka 1 sampai dengan 5. Berikut adalah penjelasan terkait proses yang ada:

1. Read Excell, bagian ini merupakan proses dimana dataset original diambil dari file excell. Proses ini memungkinkan pemilihan atribut dari dataset secara manual. Selain melakukan pemilihan secara manual, proses ini juga melakukan pemilihan tipe dari semua atribut termasuk atribut label atau atribut tujuan untuk klasifikasi.
2. Multiply, proses ini hanya memberikan lebih banyak output dari satu input yang sama. Multiply dalam penelitian ini digunakan untuk data heregistrasi mahasiswa agar dapat diakses untuk tiga jenis algoritma.
3. Bayes, dalam proses ini sebenarnya adalah proses validasi yaitu menggunakan X-Validation. Proses ini memungkinkan perulangan sebanyak 10 kali agar data yang ada dicampur secara merata. Didalam proses ini terdapat proses perhitungan algoritma naive bayes yang disertakan pula confussion matrix untuk melakukan perhitungan akurasi dari algoritma tersebut sebagaimana ada di gambar 2. Proses ini memungkinkan dataset dibagi menjadi dua bagian yaitu satu digunakan untuk data training, dan digunakan untuk data testing.
4. KNN, Dalam proses ini sebenarnya sama dengan proses pada bayes. Perbedaannya hanya ada pada bagian training yang diisi menggunakan algoritma KNN.
5. C45, Proses ini juga sama dengan proses algoritma lain seperti bayes dan KNN, proses training dalam hal ini diisi menggunakan algoritma Decission tree C45.



Gambar 2 Proses X-Validation

4.2 Validasi

Proses validasi sebenarnya telah sedikit dibahas pada tahap eksperimen. Proses validasi yang digunakan dalam penelitian ini adalah X-Validation dengan perincian menggunakan 10 folds cross validation. Proses ini memungkinkan dataset yang ada dibagi menjadi 10 bagian. Satu bagian diantaranya dijadikan data testing dengan 9 bagian lainnya menjadi data training untuk satu persatu algoritma. Proses ini berlanjut dengan menggunakan satu bagian yang lain untuk digunakan sebagai data testing. Proses ini akan berlanjut sampai

dengan iterasi yang ke 10 agar kesemua bagian data mendapatkan proporsi menjadi data testing.

4.3 Pengukuran Akurasi Algoritma

Pengukuran akurasi dalam penelitian ini menggunakan confusion matrix atau matrix kebingungan untuk semua algoritma yang dibandingkan. Aplikasi rapid miner memiliki output yang sangat mudah dibaca dengan menggunakan matrix ini. Gambar 3 merupakan hasil print screen dari aplikasi rapid miner yang menunjukkan performa algoritma KNN untuk klasifikasi heregistrasi mahasiswa. Dalam gambar 3 terlihat akurasi dari KNN sebesar 80,46%. Artinya dalam keseluruhan data testing ada 8046 yang sesuai dari 10000 percobaan record. Label dari data yang ada memiliki 2 varian yaitu Ya dan Tidak. Dalam matrix kebingungan ini terdapat tabel berwarna kuning dengan jumlah matrix sebesar 2 x 2. Matrix ini menunjukkan jumlah data testing yang sesuai dan tidak sesuai dengan label yang sebenarnya. Bagian atas merupakan bagian data atau label sebenarnya sedangkan bagian kiri merupakan bagian prediksi atau hasil perhitungan algoritma. Dalam gambar 3 true Tidak dengan pred. Tidak sebesar 25 dan true Ya dengan pred. Ya sebesar 2120. Artinya dataset hasil perhitungan algoritma yang memiliki hasil sama dengan data asli sebanyak 2120+25 yaitu 2145. Sedangkan keseluruhan jumlah record adalah 25+32+489+2120 yaitu 2666. Tingkat akurasi algoritma KNN dapat dihitung dengan 2145 dibagi 2666 dikalikan 100% yaitu 80,4576%.

Untuk algoritma naive bayes dan decision tree C45 dilakukan proses yang sama dengan proses dari KNN diatas. Gambar 4 merupakan hasil akurasi dari algoritma naive bayes dengan nilai akurasi sebesar 74,49%. Sedangkan gambar 5 merupakan hasil akurasi dari decision tree C45 dengan tingkat akurasi tertinggi dari kesemuanya yaitu 80,72%.

accuracy: 80.46% +/- 0.66% (mikro: 80.46%)			
	true Tidak	true Ya	class precision
pred. Tidak	25	32	43.86%
pred. Ya	489	2120	81.26%
class recall	4.88%	98.51%	

Gambar 3 Akurasi algoritma KNN

accuracy: 74.49% +/- 3.33% (mikro: 74.49%)			
	true Tidak	true Ya	class precision
pred. Tidak	130	296	30.52%
pred. Ya	384	1856	82.86%
class recall	25.29%	86.25%	

Gambar 4 Akurasi algoritma naive bayes

accuracy: 80.72% +/- 0.16% (mikro: 80.72%)			
	true Tidak	true Ya	class precision
pred. Tidak	0	0	0.00%
pred. Ya	514	2152	80.72%
class recall	0.00%	100.00%	

Gambar 5 Akurasi algoritma decision tree C45

5. SIMPULAN

Hasil perhitungan menunjukkan bahwa ketiga algoritma yaitu KNN, naive bayes serta decision tree C45 dapat digunakan untuk melakukan klasifikasi heregistrasi calon mahasiswa. Hasil akurasi dari penelitian ini dapat dilihat pada tabel 2 berikut.

Tabel 2 Hasil Penelitian

Algoritma	accuracy	precision	recall	AUC
KNN	80,46 %	81,26 %	98,51 %	0,578
Naive bayes	74,49 %	82,85 %	86,25 %	0,602
Decission tree C45	80,72	80,72 %	100 %	0,500

6. UCAPAN TERIMAKASIH

Penelitian ini sepenuhnya didanai oleh DRPM DIKTI dalam hibah Penelitian Dosen Pemula tahun pendanaan 2018. Ucapan terimakasih kami ucapkan kepada segenap jajaran Kementerian Riset dan Pendidikan Tinggi Republik Indonesia.

7. REFERENSI

- Alkaromi, M Adib. 2014. "Information Gain Untuk Pemilihan Fitur Pada Klasifikasi Heregistrasi Calon Mahasiswa Dengan Menggunakan K-NN."
- Amancio, D. R., C. H. Comin, D. Casanova, G. Travieso, O. M. Bruno, F. a. Rodrigues, and L. Da F. Costa. 2013. "A Systematic Comparison of Supervised Classifiers," October. <http://arxiv.org/abs/1311.0202v1>.
- Ashari, Ahmad, Iman Paryudi, and A Min Tjoa. 2013. "Performance Comparison between Naïve Bayes , Decision Tree and K-Nearest Neighbor in Searching Alternative Design in an Energy Simulation Tool" 4 (11): 33–39.
- Han, Jiawei, and Micheline Kamber. 2006. *Data Mining: Concepts and Techniques Second Edition*. Elsevier. Elsevier.
- Kusrini, Sri Hartati, Retantyo Wardoyo, and Agus Harjoko. 2009. "Perbandingan Metode Nearest Neighbor Dan Algoritma c4.5 Untuk Menganalisis Kemungkinan Pengunduran Diri Calon Mahasiswa Di Stmik Amikom Yogyakarta" 10 (1).
- Kusrini, and Luthfi Emha Taufiq. 2009. *Algoritma Data Mining*. Yogyakarta: Andi Offset.
- Larose, Daniel T. 2005. *Discovering Knowledge in Data: An Introduction to Data Mining*. John Wiley & Sons.
- Patel, Kanu, Jay Vala, and Jaymit Pandya. 2014. "Comparison of Various Classification Algorithms on Iris Datasets Using WEKA" 1 (1): 1–7.
- Pudjianto, Tacbir Hendro, Faiza Renaldi, and Age Teogunadi. 2011. "Penerapan Data Mining Untuk Menganalisa Kemungkinan Pengunduran Diri Calon Mahasiswa Baru."
- Ragab, Abdul Hamid M., Amin Y. Noaman, Abdullah S. Al-Ghamdi, and Ayman I. Madbouly. 2014. "A Comparative Analysis of Classification Algorithms for Students College Enrollment Approval Using Data Mining." *Proceedings of the 2014 Workshop on Interaction Design in Educational Environments - IDEE '14*. New York, New York, USA: ACM Press, 106–13. doi:10.1145/2643604.2643631.
- Sugianti, Devi. 2012. "Algoritma Bayesian Classification Untuk Memprediksi Heregistrasi Mahasiswa Baru Di STMIK Widya Pratama," no. 2: 1–5.
- Witten, Ian H, Eibe Frank, and Mark A. Hall. 2011. *Data Mining: Practical Machine Learning Tools and Techniques 3rd Edition*. Elsevier.
- Wu, Xindong. 2009. *The Top Ten Algorithms in Data Mining*. Edited by Vipin Kumar. New York: Taylor & Francis Group, LLC.
- Wu, Xindong, Vipin Kumar, J. Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J. McLachlan, et al. 2007. *Top 10 Algorithms in Data Mining. Knowledge and Information Systems*. Vol. 14. doi:10.1007/s10115-007-0114-2.