

PREDIKSI PENERIMA BEASISWA MENGGUNAKAN ALGORITMA C4.5 (STUDI KASUS : UNIVERSITAS PERADABAN)

Jumaroh¹, Sorikhi², Tezhar Rayendra Trastaronny Pastika Nugraha³

¹Program Studi Teknik Informatika, Fakultas Sains dan Teknologi, Universitas Peradaban,

²Program Studi Teknik Informatika, Fakultas Sains dan Teknologi, Universitas Peradaban,

³Program Studi Teknik Informatika, Fakultas Sains dan Teknologi, Universitas Peradaban,

jumaroh97@gmail.com, soryc2001@yahoo.com, tezhar.rayendra19@gmail.com

Jl. Raya Pagojengan KM 03 Paguyangan Brebes

Abstract

Kata Kunci:

*data mining,
algoritma C4.5,
beasiswa, confusion
matrix, universitas
peradaban*

Tingginya jumlah mahasiswa yang berprestasi pada perguruan tinggi dapat dimaksimalkan dengan kebijakan yang disebut beasiswa. Universitas Peradaban merupakan salah satu perguruan tinggi yang memberikan bantuan belajar berupa beasiswa salah satunya yaitu beasiswa Peningkatan Prestasi Akademik (PPA) yang terus meningkat setiap tahunnya. Melihat indikator dalam penyeleksian berkas pengajuan beasiswa dengan beberapa macam kriteria, namun masih menggunakan cara manual dalam proses penyeleksiannya, maka pemilihan mahasiswa yang berhak menerima beasiswa PPA berdasarkan kelengkapan data yang dikumpulkan terkesan objektif, kurang efektif dan membutuhkan waktu yang relatif lebih lama, untuk mengatasi permasalahan tersebut diperlukan suatu metode yang dapat memberikan keputusan yang tepat efektif dan efisien dalam proses penyeleksian dan penentuan penerima beasiswa PPA bagi seluruh pendaftar berdasarkan data yang masuk. Metode yang digunakan yaitu decision tree algoritma C4.5 dengan teknik pohon keputusan. Berdasarkan Hasil klasifikasi menggunakan algoritma C4.5 menunjukkan bahwa diperoleh akurasi mencapai 96,190476%.

Keywords

*data mining, C4.5
algorithm,
scholarship,
confusion matrix,
peradaban
university*

The high amount of students achievers at universities can be maximized with a policy called a scholarship. Peradaban University is one of the tertiary institutions that provide learning assistance in the form of students, one of which is the scholarship to increase academic achievement (PPA) which continues to increase every year. Look at the indicators in selecting the scholarship submission file with several kinds of criteria, but still using the manual method in the selection process, the selection of students who are entitled to receive PPA scholarships based on the completeness of the data collected seems objective, less effective and requires a relatively long time, to overcome problems This requires a method that can provide effective and efficient decisions in the process of selection and determination PPA scholarship recipients for all applicants based on incoming data. The method used is a decision tree C4.5 algorithm with decision tree techniques. Based on the classification results using the C4.5 algorithm shows that the accuracy is 96.190476%.

Introduction

Pendidikan merupakan salah satu faktor kemajuan dan kemandirian bangsa. Dengan majunya pendidikan suatu bangsa, maka tercipta generasi penerus yang berkualitas. Pendidikan bertujuan untuk mengembangkan potensi peserta didik agar menjadi manusia yang beriman dan bertakwa kepada Tuhan Yang Maha Esa, berakhlak mulia, sehat, berilmu, kreatif, mandiri dan menjadi warga negara yang demokratis serta bertanggung jawab [1]. Mewujudkan pendidikan yang berkualitas diperlukan biaya pendidikan yang cukup besar, oleh karena itu bagi setiap peserta didik pada setiap satuan pendidikan berhak mendapatkan biaya pendidikan bagi mereka yang orang tuanya tidak mampu membiayai pendidikannya dan berhak mendapatkan beasiswa bagi mereka yang berprestasi [2].

Tingginya jumlah mahasiswa yang berprestasi pada perguruan tinggi dapat dimaksimalkan dengan kebijakan yang disebut beasiswa [3]. Universitas Peradaban merupakan salah satu perguruan tinggi yang memberikan bantuan belajar berupa beasiswa salah satunya yaitu beasiswa Peningkatan Prestasi Akademik (PPA), data beasiswa yang ada di Universitas Peradaban tidak banyak memiliki kegunaan seolah-olah menjadi sekumpulan data terabaikan yang akan bertambah tiap tahunnya. Data tersebut hanya digunakan untuk laporan kepada pihak pengelola universitas. Data tentang beasiswa dapat memberikan informasi berguna bagi universitas jika dimanfaatkan dengan maksimal. Beasiswa PPA yang terus meningkat setiap tahunnya sedangkan kuota penerima beasiswa PPA di universitas ini terbatas. Melihat indikator dalam penyeleksian berkas pengajuan beasiswa dengan beberapa macam kriteria, namun masih menggunakan cara manual dalam proses penyeleksiannya, maka pemilihan mahasiswa yang berhak menerima beasiswa PPA berdasarkan kelengkapan

data yang dikumpulkan terkesan objektif, kurang efektif dan membutuhkan waktu yang relatif lebih lama.

Hal ini prediksi penerima beasiswa sangat penting bagi suatu perguruan tinggi. Dimana dengan adanya prediksi penerima beasiswa dapat membantu menentukan mahasiswa yang berhak atau layak untuk mendapatkan beasiswa dari perguruan tinggi. Ada banyak cara yang dapat digunakan untuk menganalisis dalam memprediksi, salah satunya yaitu data mining [3]. Data mining adalah proses pencarian pola data yang tidak diketahui atau tidak diperkirakan sebelumnya [4]. Data mining sendiri memiliki beberapa metode salah satunya yaitu klasifikasi diantaranya Naive Bayes, Decision Tree, K-NN dan Linier Regreesion. Salah satu teknik data mining yaitu decision tree yang dapat digunakan sebagai prediksi adalah algoritma C4.5 [3].

Algoritma C4.5 merupakan algoritma klasifikasi data dengan teknik pohon keputusan yang dapat mengolah data numerik dan diskrit, dapat menangani nilai atribut yang hilang, menghasilkan aturan-aturan yang mudah diinterpretasikan dan tercepat diantara algoritma-algoritma lain [5]. Kelemahan Algoritma C4.5 salah satunya terdapat di skalabilitas yaitu data training hanya dapat digunakan dan disimpan secara keseluruhan pada waktu yang bersamaan [6].

Pada studi kasus lain, berdasarkan penelitian yang dilakukan oleh Amiruddin dan Rezqiwati Ishak [7], "Prediksi Jumlah Mahasiswa Registrasi Per Semester Menggunakan Linier Regresi Pada Universitas ICHSAN Gorontalo" dimana pada penelitian tersebut dilakukan teknik prediksi menggunakan metode Linier Regresi dan MAPE. Tujuan dari penelitian ini adalah membangun aplikasi untuk memprediksi jumlah mahasiswa registrasi. Berdasarkan hasil penelitian dari 2 prodi yang dipilih yakni prodi Teknik Informatika didapatkan hasil tingkat error 4.24% atau tingkat akurasi 95.76%, dan

untuk prodi Ilmu Hukum didapatkan tingkat error 7.69% atau tingkat akurasi 92.31%, dengan demikian aplikasi yang sudah dibangun layak untuk digunakan.

Selain itu pada penelitian yang dilakukan Wiwit Supriyanti, Kusriani dan Armadiyah Amborowati [8] "Perbandingan Kinerja Algoritma C4.5 dan Naïve Bayes Untuk Ketepatan Pemilihan Kosentrasi Mahasiswa" dimana hasil uji kinerja algoritma klasifikasi untuk kasus ketepatan pemilihan konsentrasi mahasiswa untuk algoritma C4.5 tanpa penambahan seleksi fitur forward selection diperoleh nilai akurasi sebesar 84,43%, kemudian setelah ditambahkan seleksi fitur forward selection meningkat menjadi 84,98%. Sedangkan pada algoritma Naive Bayes tanpa penambahan seleksi fitur forward selection diperoleh nilai akurasi sebesar 78,47%, setelah ditambahkan seleksi fitur forward selection meningkat menjadi 82,01% dalam penelitian ini algoritma C4.5 lebih unggul dibanding Naïve Bayes.

Oleh sebab itu pada penelitian ini penulis tertarik untuk memprediksi penerima beasiswa peningkatan prestasi akademik (PPA) yang diimplementasikan pada mahasiswa Universitas Peradaban menggunakan algoritma C4.5 karena kemampuan model/rule yang dihasilkan dapat memprediksi dengan benar walaupun data ada nilai dari atribut yang hilang dan rule yang dihasilkan mudah dipahami untuk memprediksi.

Method

Metode penelitian pada penelitian ini adalah penelitian eksperimen dengan tahapan penelitian sebagai berikut :

1. Pengumpulan data

Penelitian ini diawali dengan melakukan pengumpulan data. Data yang diperoleh adalah data beasiswa di Universitas Peradaban.

2. Pengolahan data awal

Setelah data terkumpul, langkah selanjutnya adalah mengolah data yang ada. Menyeleksi atribut data yang akan

digunakan dan data yang tidak relevan dihilangkan.

3. Model/ metode yang diusulkan

Peneliti mengusulkan model yang akan digunakan. Model yang digunakan tersebut berupa metode dalam teknik prediksi data mining yaitu klasifikasi dengan algoritma C4.5 yang dioptimisasi untuk pemilihan atribut yang digunakan.

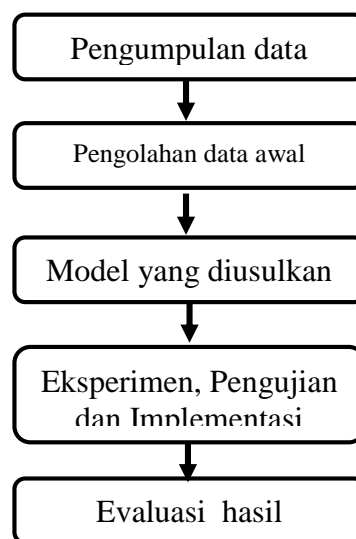
4. Eksperimen, pengujian model dan implementasi

Pengujian model, dari algoritma yang sudah ditentukan maka dataset yang ada diolah. Sehingga menghasilkan model yang diinginkan. Model yang dihasilkan selanjutnya diimplementasi dalam program.

5. Evaluasi dan validasi hasil

Setelah dilakukan eksperimen terhadap semua dataset dengan model yang diusulkan maka akan menghasilkan nilai-nilai akurasi dan performa. Kemudian hasil tersebut dianalisa dan dievaluasi menggunakan confusion matrix. Dari hasil evaluasi dapat ditarik kesimpulan dari penelitian dan eksperimen ini

Pada penelitian ini akan mengusulkan sebuah metode klasifikasi data untuk prediksi penerima beasiswa. Hasil dataset yang baru akan diolah menggunakan algoritma C4.5 dan akan dibangun pohon keputusan untuk menghasilkan sebuah rule.



Gambar 1. Metode Penelitian

Results and Discussion

Pada penelitian ini telah dilakukan pengujian dengan menggunakan *python* pada data beasiswa PPA Universitas Peradaban sebanyak 105 orang, dimana atribut yang didapat nama mahasiswa, status, nama mahasiswa sesuai penulisan pada rekening, nama bank, jenis kelamin, nim, jurusan, S1/D4/D3, semester, IPK, kedisiplinan, keaktifan. Pada penelitian ini dilakukan klasifikasi guna untuk melakukan perhitungan. Berikut klasifikasi IPK dapat dilihat Tabel 1.

Tabel 1 Tabel Klasifikasi IPK

Klasifikasi IPK	
≥ 3.50	Sangat Memuaskan
≥3.10	Memuaskan
<3.09	Kurang

Pada klasifikasi selanjutnya yaitu klasifikasi kedisiplinan yang dapat dilihat pada Tabel 2 Tabel klasifikasi kedisiplinan.

Tabel 2 Tabel Klasifikasi Kedisiplinan

Klasifikasi Kedisiplinan	
V	Disiplin
K	Kurang Disiplin
X	Tidak Disiplin

Pada klasifikasi selanjutnya yaitu klasifikasi keaktifan yang dapat dilihat pada Tabel 3 Tabel klasifikasi keaktifan.

Tabel 3 Tabel Klasifikasi Keaktifan

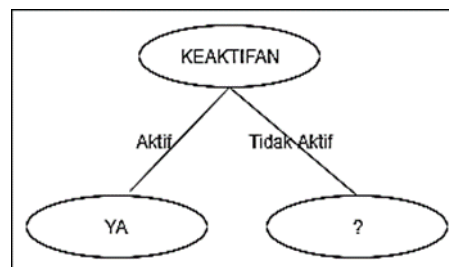
Klasifikasi Keaktifan	
V	Aktif
X	Tidak Aktif

Selanjutnya menghitung *entropy* total, *entropy* masing - masing atribut dan menghitung gain dan menentukan gain tertinggi, hal ini dapat dilihat dari Tabel node 1 berikut ini :

Tabel 4 Tabel Perhitungan Node 1.

NO-DE 1	ATRIBUT	KELAS	JUMLAH KASUS	LAYAK	TIDAK LAYAK	ENTROPY	GAIN
	TOTAL		105	86	19	0,682156	
	IPK						0,230622
		Sangat Memuaskan	51	51	0	0	
		Memuaskan	52	35	17	0,911752	
		Kurang	2	0	2	0	
	KEDISIPLINAN						0,485961
		Disiplin	78	78	0	0	
		Kurang Disiplin	17	13	4	0,787127	
		Tidak Disiplin	10	8	2	0,721928	
	KEAKTIFAN						0,513405
		Aktif	83	83	0	0	
		Tidak Aktif	24	5	19	0,738285	

Dari hasil pada Tabel 4 menghasilkan pohon keputusan yang terbentuk dari pencarian *node* 1 adalah seperti pada Gambar 2.



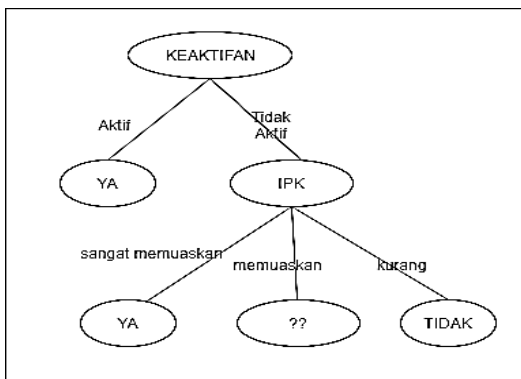
Gambar 2 Gambar Pohon Keputusan Hasil Perhitungan Node 1

Kemudian data di Tabel 4 dihitung lagi *entropy* atribut keaktifan - tidak aktif dan *entropy* setiap atribut serta *gain* ya, sehingga hasilnya seperti pada Tabel 5.

Tabel 5 Tabel Perhitungan Node 1.1

NODE	ATRIBUT	KELAS	JUMLAH KASUS	LAYAK	TIDAK LAYAK	ENTROPY	GAIN
1.1	Keaktifan	Tidak Aktif	24	5	19	0,738285	
	IPK						0,506127
		Sangat Memuaskan	4	4	0	0	
		Memuaskan	18	1	17	0,309543	
		Kurang	2	0	2	0	
	Kedisiplinan						0,24835
		Disiplin	12	5	7	0,979869	
		Kurang	4	0	4	0	
		Tidak Disiplin	8	0	8	0	

pohon keputusan yang terbentuk dari pencarian *node* 1.1 adalah seperti pada Gambar 3.



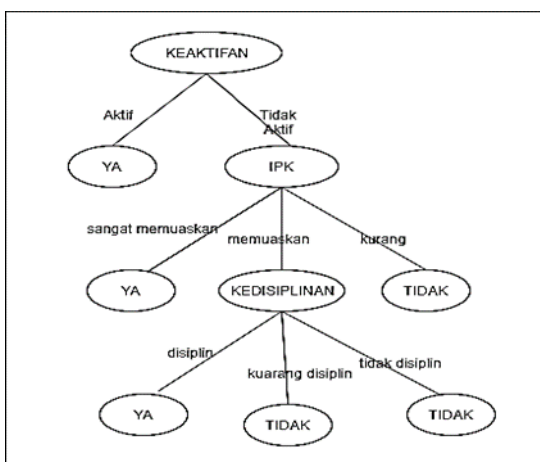
Gambar 3 Gambar Pohon Keputusan Hasil Perhitungan Node 1.1

Kemudian data di Tabel 5 dihitung lagi *entropy* atribut IPK - memuaskan dan *entropy* setiap atribut serta *gain* ya, sehingga hasilnya seperti pada Tabel 6.

Tabel 6 Tabel Perhitungan Node 1.1.1

NODE	ATRIBUT	KELAS	JUMLAH KASUS	LAYAK	TIDAK LAYAK	ENTROPY	GAIN
1.1.1	IPK	Memuaskan	18	1	17	0.309543	
	KEDISIPLINAN						0.092869
		Disiplin	6	1	5	0.650022	
		Kurang Disiplin	4	0	4	0	
		Tidak Disiplin	7	0	7	0	

Adapun pohon keputusan yang terbentuk dari pencarian node 1.1.1 adalah seperti pada Gambar 4.



Gambar 4 Gambar Pohon Keputusan Hasil Perhitungan Node 1.1.1

Pada Gambar 4 adalah perhitungan *gain* dan *entropy* setiap variabel sampai

menghasilkan *node* 1.1.1 adalah sebagai berikut :

1. Jika aktif maka ya
2. Jika tidak aktif dan IPK sangat memuaskan maka ya
3. Jika tidak aktif dan IPK kurang maka tidak
4. Jika IPK memuaskan maka dan disiplin maka ya
5. Jika kurang disiplin maka tidak
6. Jika tidak disiplin maka tidak

Pada tahap selanjutnya data akan di proses menggunakan python dengan IDE jupyter notebook, maka data akan dilihat nilai akurasi dan rule nya data import di python dengan menggunakan format csv. Berikut Tabel 7 adalah tabel data yang akan di import ke dalam python :

Tabel 4. 10 Tabel Data yang diimport ke python.

IPK	Kedisiplinan	Keaktifan	Decison
3.67	V	V	Yes
3.86	K	V	Yes
3.32	V	V	Yes
3.42	V	V	Yes
.....
3.32	X	X	No
3.56	X	X	No
3.47	X	X	No

Keterangan pada kolom kedisiplinan :

V : disiplin

K : kurang disiplin

X : tidak disiplin

Keterangan pada kolom keaktifan :

V : aktif

X : tidak aktif

Berikut adalah proses pengolahan data menggunakan algoritma C4.5 untuk mendapatkan pohon keputusan atau rule dengan python. Langkah pertama yaitu menyiapkan framework Chefboost. Chefboost adalah kerangka pohon pengambilan keputusan yang diaktifkan dengan gradient boosting, random forest dan adaboost yang diaktifkan termasuk algoritma ID3, C4.5, CART dan regression tree algorithms dengan dukungan fitur kategorikal. Berikut adalah coding untuk mendapat rule Algoritma C4.5.

```
import Chefboost as chef
import pandas as pd
```

selanjutnya mengimportkan dataset dengan coding berikut :

```
df = pd.read_csv("dataset/DataahLabel.csv")
df.head()
```

selanjutnya dengan coding berikut untuk mendapatkan *rule* dan akurasi algoritma C4.5

```
config = {'algorithm': 'C4.5'}
model = chef.fit(df.copy(), config)
```

Berikut adalah *rule* yang diperoleh dari *python*

```
def findDecision(obj): #obj[0]: IPK, obj[1]: DISIPLIN, obj[2]: AKTIF
    if obj[2] == 'V':
        return 'yes'
    elif obj[2] == 'X':
        if obj[1] == 'V':
            if obj[0] > 3.09:
                return 'yes'
            elif obj[0] <= 3.09:
                return 'no'
        elif obj[1] == 'X':
            return 'no'
        elif obj[1] == 'K':
            return 'no'
```

Rule diatas yang nantinya akan digunakan pada implementasi program. Dari *rule* yang diperoleh dari *python*, adapun *rule* atau aturan yang terbentuk adalah sebagai berikut :

1. Jika aktif maka ya
2. Jika tidak aktif, disiplin dan IPK >3.09 maka ya
3. Jika IPK <3.09 maka tidak
4. Jika kurang disiplin maka tidak
5. Jika tidak disiplin maka tidak

Selanjutnya pengujian ini dilakukan untuk mengetahui kinerja dari algoritma C4.5 dalam melakukan klasifikasi terhadap kelas yang ditentukan dalam penelitian ini. Pengujian dilakukan menggunakan *confusion matrix*.

Berikut adalah akurasi yang diperoleh dari *python*

```
C4.5 tree is going to be built...
Accuracy: 96.19047619047619 % on 105 instance
finished in 2.2079503536224365 seconds
```

Algoritma C4.5 melakukan training terhadap data-data yang telah dibagi (*split validation*). Data yang diimportkan sebanyak 105 record. Berikut tabel perhitungan

confusion matrix yang diperoleh ditunjukkan pada Tabel 8.

Tabel 8 Tabel perhitungan *confusion matrix*

Correct Classification	Classified negative (yes)	Classified positive (no)
Actual negative (yes)	15	4
Actual positive (no)	0	86

Akurasi yang dihasilkan dari pengujian pada *python* sebesar 96,190476% Berikut adalah perhitungan akurasi :

$$A = \left(\frac{a+d}{a+b+c+d} \right) \times 100\% = \dots\dots\%$$

$$A = \left(\frac{15+86}{15+4+0+86} \right) \times 100\% = 96,190476\%$$

Nilai *precision* dihitung dengan cara membagi jumlah data benar yang bernilai positif (*True Positive*) dibagi dengan jumlah data benar yang bernilai positif (*True Positive*) dan data salah yang bernilai negatif (*False Negative*).

Berikut perhitungan *precision* :

$$P = \left(\frac{d}{b+d} \right) \times 100\% = \dots\dots\%$$

$$P = \left(\frac{86}{4+86} \right) \times 100\% = 95,55\%$$

Nilai *recall* dihitung dengan cara membagi jumlah data benar yang bernilai positif (*True Positive*) dibagi dengan jumlah data benar yang bernilai positif (*True Positive*) dan data salah yang bernilai positif (*False Positive*)

$$R = \left(\frac{d}{c+d} \right) \times 100\% = \dots\dots\%$$

$$R = \left(\frac{86}{0+86} \right) \times 100\% = 100\%$$

Pengujian menghasikan akurasi, *precision*, dan *recall* yang dapat dilihat pada Tabel 9 :

Tabel 9 Hasil Pengujian

Akurasi	<i>precision</i>	<i>Recall</i>
96,190476%	95,55%	100%

Berdasarkan tabel yang telah dijelaskan di atas, maka dapat diketahui bahwa pada *confusion matrix* memiliki nilai akurasi 96,190476%, *Precision* 95,55% dan *recall* 100%. Berdasarkan proses pengolahan data menggunakan Algoritma C4.5 yang telah menghasilkan pohon keputusan dan *rule* yang telah terbentuk, selanjutnya akan diimplementasikan *rule* tersebut untuk

membuat program prediksi penerima beasiswa berbasis web.

Conclusion and Suggestions

Berdasarkan hasil penelitian yang telah dilakukan oleh peneliti, maka dapat disimpulkan bahwa penerima beasiswa dapat diprediksi dan dievaluasi dengan memanfaatkan teknik data mining menggunakan algoritma *decision tree* C4.5 untuk memprediksi (menentukan kelas) dari data *training* yang telah diperoleh. Hasil percobaan dan pengujian prediksi penerima beasiswa dengan *python* menggunakan metode *decision tree* C4.5, diperoleh akurasi sebesar 96,190476% dengan kriteria akurasi *excellent classification* menggunakan *confusion matrix*.

Saran yang dapat disampaikan untuk meningkatkan kinerja dan menyempurnakan penelitian yang telah dibuat, peneliti memberikan saran sebagai berikut :

1. Penelitian ini dapat dikembangkan dengan menggabungkan atau membandingkan dengan algoritma klasifikasi lain seperti K-Nearest Neighbors, Support Vector Machines, Naive Bayes dan yang lainnya untuk mendapatkan hasil prediksi yang lebih baik.
2. Pada penelitian lebih lanjut disarankan untuk mencoba algoritma lain seperti Predictive modelling algoritmanya Linear Regression, Neural Network, Support Vector Machine, dan lain-lain..

References

- [1] R. Gunawan, "Implementasi Data Mining Untuk Memprediksi Prestasi Siswa Berdasarkan Status Sosial Dan Kedisiplinan Pada Smk Bayu Pertiwi Menggunakan Metode Regresi Linier Berganda," *Sains dan Komputer (SAINTIKOM)*, pp. 175-183, 2018.
- [2] Rismayanti, "IMPLEMENTASI ALGORITMA C4.5 UNTUK MENENTUKAN PENERIMA BEASISWA

DI STT HARAPAN MEDAN," *Jurnal Media Infotama*, vol. 12, p. 2, 2016.

- [3] Rahman, Muhammad Arif, "ALGORITMA C45 UNTUK MENENTUKAN MAHASISWA PENERIMA BEASISWA (STUDI KASUS : PPS IAIN RADEN INTAN BANDAR LAMPUNG)," *Jurnal TIM Darmajaya*, vol. 1, p. 2, 2015.
- [4] E. R. Paramita Mayadewi, "Prediksi Nilai Akhir Proyek Mahasiswa Menggunakan Algoritma Klasifikasi Data Mining," *Seminar Nasional Sistem Informasi Indonesia*, 2015.
- [5] Rizky Haqmanullah Pambudi, Budi Darma Setiawan, Indriati, "Penerapan Algoritma C4.5 Untuk Memprediksi Nilai Kelulusan Siswa Sekolah Menengah Berdasarkan Faktor Eksternal," *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, vol. 2, pp. 2637-2643, 2018.
- [6] Yuni Sara Luvia, Dedy Hartama, Agus Perdana Windarto, Solikhun, "Penerapan Aloritma C4.5 Untuk KLasifikasi Predikat Keberhasilan Mahasiswa Di AMIK TUNAS BANGSA," *JURASIK (Jurnal Riset Sistem Informasi & Teknik Informatika)*, vol. 1, 2016.
- [7] Amiruddin, Rezqiwati Ishak, "Prediksi Jumlah Mahasiswa Registrasi Per Semester Menggunakan Linier Regresi Pada Universitas ICHSAN Gorontalo," *ILKOM Jurnal Ilmiah*, vol. 10, 2 Agustus 2018.
- [8] Wiwit Supriyanti, Kusri dan Armadyah Amborowati, "Perbandingan Kinerja Algoritma C4.5 dan Naive Bayes Untuk Ketepatan Pemilihan Kosentrasi Mahasiswa," *Jurnal INFORMA Politeknik Indonusa Surakarta*, Vols. 2442-7942, p. 3, 2016.