

KARAKTERISTIK SISWA SISWI SMA YANG MENENTUKAN PEMILIHAN PERGURUAN TINGGI SWASTA DI INDONESIA DENGAN MENGGUNAKAN LOGISTIK REGRESI

Hendra Achmadi^{1*}, Indah Fatmawati²⁾, Sylvia Samuel³⁾

Universitas Pelita Harapan, Banten, Indonesia

Hendra.achmadi@uph.edu

infatmawati@yahoo.com

Sylvia.Samuel@uph.edu

ABSTRACT

Jumlah Calon mahasiswa yang masuk pada perguruan tinggi khususnya perguruan tinggi swasta adalah merupakan masalah serius. Dan jumlah penurunan calon mahasiswa yang terjadi dalam masa pandemic covid-19 antara tahun 2018-2019 dan berlanjut sampai tahun 2022 adalah merupakan masalah serius bagi perguruan tinggi Swasta di Indonesia. Karena itu dalam penelitian ini menfokuskan pada mencari karakteristik utama dari siswa siswi SMA dalam melakukan pemilihan perguruan tinggi swasta di Jakarta dan Sekitarnya. Metode penelitian yang dipakai adalah data mining, dan data yang dipakai adalah berasal dari data primer yang didapat dari kuesioner yang dibagikan kepada siswa-siswi SMA kelas 11 dan 12 di Jakarta dan sekitarnya, dan didapat 438 responden dan setelah dilakukan data cleansing didapat 295 responden. Dengan menggunakan metode Logistic Regression dan Supervisor Learning akan didapat model dari keputusan siswa siswi SMA dalam memilih perguruan tinggi swasta. Hasil dari penelitian ini adalah Dari hasil didapat bahwa keputusan pemilihan Universitas, yang dianalisa dari nilai Beta yang positif dari Logistik Regression. Dari hasil didapat bahwa karakteristik pertama adalah Tipesks atau Tipe sekolah (Swasta atau Negeri) dan memang yang banyak masuk dalam Universitas Swasta dari sekolah swasta. Karakteristik kedua adalah UANGSAKU, disini menggambarkan akan taraf hidup dari keluarga, dalam penelitian ini UANGSAKU yang paling banyak adalah diantara < 1 juta sebanyak 237 responden, Berikutnya adalah banyaknya presentasi dari pihak universitas ke sekolah juga sangat berpengaruh terhadap siswa siswi di SMA, kemudian karakteristik selanjutnya adalah SEX yang didominasi oleh PRIA sebanyak 151 orang dari 295 responden. Dilanjutkan dengan DOMISILI yang didominasi oleh daerah Tangerang, Banten, Bekasi dan Bogor sebanyak 144 responden, dan yang terakhir adalah TRANSPORT yang didominasi oleh antar jemput orang tua dan grab.

Kata Kunci: Karakteristik siswa-siswi SMA, Data Mining, Logistic Regression

1. PENDAHULUAN

Dalam tiga tahun terakhir, industri jasa pendidikan di Indonesia telah berkembang dengan pesat. Sebuah bukti nyata adalah adanya peningkatan sebesar 42,55% dari jumlah lembaga pendidikan di Indonesia antara tahun 2017 dan 2018. Di Provinsi Banten, ada peningkatan sebesar 38,65% pada tahun 2017 dan 2018. Peningkatan ini disebabkan oleh penambahan jenis lembaga pendidikan seperti politeknik, yang memberikan lebih banyak pilihan bagi lulusan SMA. Oleh karena itu, lembaga pendidikan tinggi harus kompetitif dalam menarik mahasiswa baru dari lulusan SMA, hal ini masih berlanjut sampai tahun 2022.

Kompetisi di antara institusi pendidikan tinggi akan mendorong upaya pemasaran yang lebih efektif untuk menjangkau siswa SMA dan membujuk mereka untuk memilih universitas atau politeknik sebagai destinasi pendidikan mereka setelah lulus SMA. Di Indonesia, institusi pendidikan tinggi dapat dibagi menjadi universitas dan politeknik, sehingga lulusan SMA memiliki banyak pilihan dalam melanjutkan Pendidikan.

Pertumbuhan universitas di seluruh Indonesia mengalami perubahan selama tiga tahun terakhir. Pada periode 2017 hingga 2018, terjadi peningkatan sebesar 42%, namun pada 2018 hingga 2019 terjadi penurunan sebesar 1%. Sama halnya dengan Provinsi Banten, rata-rata jumlah mahasiswa baru per lembaga di sana juga mengalami penurunan pada tahun 2018 hingga 2019. Peningkatan jumlah universitas di Provinsi Banten terjadi karena beberapa universitas dari Jakarta pindah ke area banten atau mendirikan cabang baru di daerah tersebut. Hal ini menyebabkan persaingan antar universitas semakin ketat. Kementerian Riset Teknologi dan Pendidikan Tinggi tahun 2018 menunjukkan terjadinya penurunan jumlah mahasiswa baru di Provinsi DKI Jakarta pada tahun 2019 sebesar 40%, dan sampai tahun 2022 sudah ada kenaikan tetapi masih terdapat penurunan di beberapa perguruan tinggi sebesar 22 %.

Jumlah mahasiswa baru adalah merupakan hal yang sangat penting bagi perguruan tinggi, agar tetap bisa mempertahankan eksistensinya di industri Pendidikan Tinggi. Calon mahasiswa di perguruan tinggi adalah siswa siswi SMA yang berusia 17 dan 18 tahun, dan termasuk Generasi Z. Menurut Kusumaningtyas et al., (2020), mengatakan bahwa generasi Z baik dalam literasi teknologi. Selain itu, penelitian lain menyebutkan bahwa gaya generasi Z memiliki kekhasan dalam pengambilan keputusan (Santoso & Triwijayati, 2018). Salah satunya adalah *online windows shopping*. Oleh karena itu dibutuhkan satu cara untuk dapat menentukan karakteristik siswa siswa yang dapat diambil untuk dapat membantu marketing di perguruan tinggi untuk dapat menambah jumlah calon mahasiswa baru di perguruan tinggi swasta di masa pandemic covid-19. Cara ini diyakini dapat juga dipakai setelah pandemi berakhir.

2. LANDASAN TEORI

2.1 Langkah Proses Data Mining

Proses Data Mining dilakukan dengan persiapan data dan dilanjutkan dengan data pemrosesan atau pembersihan data, di sini persiapan data dimulai untuk diproses lebih lanjut, untuk contoh apakah data memiliki jenis nomor atau faktor atau tanggal, dan kemudian data dalam data pembersihan juga dilakukan dengan menghilangkan karakter khusus, kemudian setelah itu dilakukan transformasi dilakukan yaitu mengubah data dari cleansing data menjadi data target yaitu proses selanjutnya adalah melakukan data mining atau model data berdasarkan metode yang cocok untuk data tersebut, dan yang terakhir adalah proses interpretasi pengetahuan yang diperoleh dari pengolahan data (Jiawei, 2012).

2.2. Logistik Regresi

Logistik Regresi menurut Janzen & Stern, (1998) dapat digunakan untuk melakukan Analisa multivaritat analisis. Dan juga menurut Fang, (2013) menyatakan bahwa logistic

regression lebih baik dibandingkan dengan linear regression, karena dalam pembuatan model logistic regression sudah ada keputusan yang harus diambil. Sedangkan menurut Sperandea et al., (2022) menyatakan bahwa logistic regression dapat dipakai dalam melakukan simulasi pendapat dari responden.

Setelah didapat karakteristik yang utama dari Logistik regresi, maka Langkah selanjutnya dapat dilanjutkan dengan membuat pohon keputusan. Selanjutnya menurut Charbuty & Abdulazeez, (2021) Algoritme pohon keputusan adalah algoritma pembelajaran mesin terawasi yang digunakan untuk tugas klasifikasi dan regresi. Ini menciptakan model keputusan seperti pohon dan kemungkinan konsekuensinya berdasarkan serangkaian fitur masukan. Algoritme secara iteratif mempartisi data menjadi subset yang lebih kecil dan lebih kecil berdasarkan nilai fitur, membuat struktur pohon di mana setiap node internal mewakili aturan keputusan berdasarkan nilai fitur dan setiap node daun mewakili label kelas atau nilai prediksi.

Algoritma pohon keputusan bekerja dengan mempartisi dataset secara rekursif menjadi subset berdasarkan nilai fitur input (Duarte et al., 2022). Pada setiap langkah proses partisi, algoritme memilih fitur yang paling baik memisahkan dataset ke dalam kelas variabel target atau meminimalkan varian variabel target. Fitur ini kemudian digunakan untuk membuat simpul di pohon keputusan. Algoritme terus mempartisi dataset di setiap node hingga kriteria penghentian terpenuhi, seperti ketika semua instance dalam subset milik kelas variabel target yang sama.

Gkikas et al., (2022) menyatakan salah satu keuntungan dari algoritma pohon keputusan adalah kemampuan interpretasinya, karena jalur keputusan dapat dengan mudah divisualisasikan dan dipahami. Namun, pohon keputusan rentan terhadap overfitting, terutama bila pohonnya dalam dan kompleks. Untuk menghindari overfitting, teknik seperti pemangkasan, regularisasi, dan metode ansambel seperti hutan acak dan peningkatan gradien dapat diterapkan.

2.3 Supervisor Learning

Lebih jauh lagi peneliti lainnya menyatakan bahwa teknik klasifikasi biasanya adalah program komputer yang belajar dari data input yang diberikan, dan menggunakan data pelatihan ini dengan tujuan untuk belajar mengklasifikasikan berdasarkan pola pengamatan pada data tersebut (Duarte et al., 2022). Di sisi lain pembelajaran terawasi untuk regresi adalah seperangkat algoritma yang digunakan untuk memprediksi nilai kontinu.

Menurut Charbuty & Abdulazeez, (2021), Algoritma pohon keputusan adalah bagian dari keluarga algoritma pembelajaran yang diawasi, dan tujuan utamanya adalah untuk membangun model pelatihan yang dapat digunakan untuk memprediksi kelas atau nilai variabel target melalui aturan keputusan pembelajaran yang disimpulkan. dari data pelatihan.

Menurut Müller & Guido, (2017) Pembelajaran yang diawasi adalah jenis pembelajaran mesin di mana algoritme belajar dari kumpulan data berlabel untuk membuat prediksi atau keputusan tentang data baru yang tidak terlihat. Dalam pembelajaran terawasi, algoritme dilatih pada sekumpulan data input dan data output yang sesuai, juga dikenal

sebagai label. Algoritme belajar untuk memetakan data input ke data output dengan menggeneralisasi pola dalam data pelatihan. Pembelajaran terbimbing menjadi area untuk banyak aktivitas penelitian dalam pembelajaran mesin. Banyak dari teknik pembelajaran yang diawasi telah menemukan penerapannya dalam pemrosesan dan analisis berbagai data

3. METODOLOGI

Penelitian ini bersifat kuantitatif, yang pertama adalah untuk mengetahui gambaran dari setiap profil pelanggan yang tertunda yang diambil melalui kuesioner kepada 202 responden dengan menggunakan *Google form*, kemudian dilakukan pengolahan data dan pembersihan data dengan menggunakan metode data mining, sehingga dapat diketahui dari lima belas karakteristik atau fitur, dimana dari fitur tersebut penting untuk menentukan keputusan, dengan menggunakan algoritma *random forest*, akan menggunakan algoritma pohon keputusan. Untuk membuat perhitungan algoritma *random forest* menggunakan program *python*, dan untuk membuat algoritma pohon keputusan menggunakan *python* juga.

4. HASIL

Persiapan Data

Data diambil data primer dari questioner yang dibagikan kepada para siswa siswa SMA kelas 11 dan 12 di daerah Jakarta dan sekitarnya dengan menggunakan google form, dan didapat 438 responden , dan kemudian dilakukan data cleansing dan tersisa 295 responden

	SMA	TIPESKS	SEX	GRADE	JURUSAN	DOMISILI	UANGSAKU	TRANSPORT	BIMBEL	PEKERJAAN	DIDIKPP	DIDIKIBU	SOCIALMEDIA	PRESENTASI	UNIV1
0	UPH College	SMS SWASTA	PRIA	XI	IPA	DKI Jakarta	1-3 Juta	Antar jemput dengan sopir	Ya	Wiraswasta	S2	S1	Ya	Ya	UPH
1	UPH College	SMS SWASTA	WANITA	XII	IPA	DKI Jakarta	1-3 Juta	Mobil	Ya	Dosen/Guru	S1	S1	Ya	Ya	UPH
2	UPH College	SMA NEGERI	WANITA	XII	IPS	DKI Jakarta	1-3 Juta	Antar jemput dengan sopir	Ya	Dosen/Guru	S1	S1	Ya	Ya	ITB
3	SMA 2	SMS SWASTA	WANITA	XII	IPA	Tangerang, Banten, Bekasi, Bogor (Termasuk Ja...	1-3 Juta	Mobil	Tidak	Wiraswasta	S3	S1	Ya	Ya	UPH
4	SMA 2	SMA SWASTA	PRIA	XII	IPA	Tangerang, Banten, Bekasi, Bogor (Termasuk Ja...	< 1 Juta	Sepeda	Tidak	Wiraswasta	SMA	SMA	Ya	Ya	UPH
...
289	SMAI PB SOEDIRMAN 1 BEKASI	SMA NEGERI	PRIA	XI	IPA	Tangerang, Banten, Bekasi, Bogor (Termasuk Ja...	< 1 Juta	Naik Kendaraan Umum	Tidak	Karyawan Swasta	SMA	SMA	Ya	Ya	UI
290	SMAK 5 Penabur	SMA SWASTA	WANITA	XII	IPA	Jawa	1-3 Juta	Mobil	Ya	Wiraswasta	S1	SMA	Ya	Ya	PRASMUL
291	SMAK 5 Penabur	SMA NEGERI	PRIA	XI	IPA	Tangerang, Banten, Bekasi, Bogor (Termasuk Ja...	< 1 Juta	Naik Kendaraan Umum	Tidak	Karyawan Swasta	SD	sederajat (SMP)	Ya	Ya	UGM
292	SMAI PB SOEDIRMAN 1 BEKASI	SMA NEGERI	WANITA	XII	IPS	Tangerang, Banten, Bekasi, Bogor (Termasuk Ja...	< 1 Juta	Naik Kendaraan Umum	Tidak	PNS	SMP	SMA	Ya	Ya	PTS
293	SMA Kolese Kanisius	SMA NEGERI	PRIA	XII	IPA	Tangerang, Banten, Bekasi, Bogor (Termasuk Ja...	< 1 Juta	Motor	Ya	Wiraswasta	S1	S1	Ya	Ya	IPB

Gambar 1. Data Sumber dari Data Primer

Encoding

Data yang sudah didapat kemudian dilakukan proses encoding, atau merubah dari tipe data kategorikal ke numerik, supaya dapat diproses selanjutnya

```

from sklearn.preprocessing import LabelEncoder
enc= LabelEncoder()
data['SMA']=enc.fit_transform(data['SMA'].values)
data['TIPESKS']=enc.fit_transform(data['TIPESKS'].values)
data['SEX']=enc.fit_transform(data['SEX'].values)
data['GRADE']=enc.fit_transform(data['GRADE'].values)
data['JURUSAN']=enc.fit_transform(data['JURUSAN'].values)
data['DOMISILI']=enc.fit_transform(data['DOMISILI'].values)
data['UANGSAKU']=enc.fit_transform(data['UANGSAKU'].values)
data['TRANSPORT']=enc.fit_transform(data['TRANSPORT'].values)
data['BIMBEL']=enc.fit_transform(data['BIMBEL'].values)
data['PEKERJAAN']=enc.fit_transform(data['PEKERJAAN'].values)
data['DIDIKPP']=enc.fit_transform(data['DIDIKPP'].values)
data['DIDIKIBU']=enc.fit_transform(data['DIDIKIBU'].values)
data['SOCIALMEDIA']=enc.fit_transform(data['SOCIALMEDIA'].values)
data['PRESENTASI']=enc.fit_transform(data['PRESENTASI'].values)
data['UNIV1']=enc.fit_transform(data['UNIV1'].values)
    
```

Gambar 2. Encoding Proses

Hasil dari proses encoding dapat di lihat di gambar 3. Hasil proses encoding

```

df1 = pd.read_csv('/SPENDING/HASILKARI1.CSV', delimiter=';')
df1
    
```

[134] ✓ 0.0s

	SMA	TIPESKS	SEX	GRADE	JURUSAN	DOMISILI	UANGSAKU	TRANSPORT	BIMBEL	PEKERJAAN	DIDIKPP	DIDIKIBU	SOCIALMEDIA	PRESENTASI	UNIV1	DIS
0	55	3	0	1	3	0	0	4	2	32	6	4	1	1	42	0
1	55	3	2	2	3	0	0	20	2	7	5	4	1	1	42	0
2	55	1	2	2	4	0	0	4	2	7	5	4	1	1	10	0
3	8	3	2	2	3	4	0	20	1	32	7	4	1	1	42	0
4	8	2	0	2	3	4	2	35	1	32	10	7	1	1	42	0
...
289	32	1	0	1	3	4	2	26	1	13	10	7	1	1	31	0
290	36	2	2	2	3	2	0	20	2	32	5	7	1	1	20	0
291	36	1	0	1	3	4	2	26	1	13	8	13	1	1	30	0
292	32	1	2	2	4	4	2	26	1	16	12	7	1	1	21	0
293	16	1	0	2	3	4	2	21	2	32	5	4	1	1	9	0

294 rows × 16 columns

Gambar 3. Hasil Proses Encoding

```

y=df1['DIS']
y
    
```

[136] ✓ 0.0s

```

... 0 0
     1 0
     2 0
     3 0
     4 0
     ..
    289 0
    290 0
    291 0
    292 0
    293 0
    Name: DIS, Length: 294, dtype: int64
    
```



```

import sklearn.metrics as met
[146] ✓ 0.0s

confusionmatrix=met.confusion_matrix(y_test,y_prediksi)
print(confusionmatrix)
[147] ✓ 0.0s
... [[42  0]
     [17  0]]

score= model1.score(X_test,y_test)
print (score)
[149] ✓ 0.0s
... 0.711864406779661
    
```

Gambar 6. Angka

Dari gambar 6 terhitung akurasi dari model sebesar 71.1 % , dimana model ini dapat memprediksi karakteristik yang akan menentukan pemilihan perguruan tinggi oleh siswa-siswi SMA sebesar 71,1 %.

KARAKTERISTIK	SKOR
TIPESKS	0,3573
UANGSAKU	0,2273
PRESENTASI	0,2079
SEX	0,0492
DOMISILI	0,0407
TRANSPORT	0,0032
SMA	-0,0112
PEKERJAAN	-0,0120
BIMBEL	-0,0286
DIDIKIBU	-0,0359
DIDIKPP	-0,0507
SOCIALMEDIA	-0,1090
JURUSAN	-0,1852
GRADE	-0,2761

Gambar 7. Beta dari Karakteristik siswa-siswi SMA

Dari gambar 7. Didapat koefisien beta dari karakteristik siswa siswa SMA yang melakukan pemilihan perguruan tinggi swasta di Jakarta dan sekitarnya.

5. HASIL

Dari logistic Regression maka akan dibuat model yaitu $Y = -0.1761 + 0.3573TIPESKS + 0.2273UANGSAKU + 0.2079PRESENTASI + 0.0492SEX + 0.0407DOMISILI + 0.0032TRANSPORT + -0.0112SMA - 0.0120PEKERJAAN - 0.0286BIMBEL - 0.0359DIDIKIBU - 0.0507DIDIKPP - 0.1090SOCIALMEDIA - 0.1852JURUSAN - 0.2761GRADE$

Dari hasil dapat dilihat bahwa Tipesks menduduki peringkat pertama, dan UANGSAKU di peringkat ke 2 dan SEX di peringkat ke 3 dan DOMISILI di peringkat ke 4 dan TRANSPORT di peringkat ke 5 yang akan berpengaruh terhadap keputusan dalam pemilihan universitas

Dari gambar 6 terhitung akurasi dari model sebesar 71.1 % , dimana model ini dapat memprediksi karakteristik yang akan menentukan pemilihan perguruan tinggi oleh siswa-siswi SMA sebesar 71,1 %.

6. PEMBAHASAN

Dari hasil didapat bahwa keputusan pemilihan Universitas , yang dianalisa dari nilai Beta yang positif dari Logistik Regresi. Dari hasil didapat bahwa karakteristik pertama adalah Tipesks atau Tipe sekolah (Swasta atau Negeri) dan memang yang banyak masuk dalam Universitas Swasta dari sekolah swasta. Karakteristik kedua adalah UANGSAKU, disini menggambarkan akan taraf hidup dari keluarga, dalam penelitian ini UANGSAKU yang paling banyak adalah diantara < 1 juta sebanyak 237 responden, Berikut nya adalah banyaknya presentasi dari pihak universitas ke sekolah juga sangat berpengaruh terhadap siswa siswi di SMA, kemudian karakteristik selanjutnya adalah SEX yang didominasi oleh PRIA sebanyak 151 orang dari 295 responden. Dilanjutkan dengan DOMISILI yang didominasi oleh daerah Tangerang, Banten, Bekasi dan Bogor sebanyak 144 respoden, dan yang terakhir adalah TRANSPORT yang dinominasi oleh antar jemput orang tua dan grab.

REFERENCES

- Charbuty, B., & Abdulazeez, A. (2021). Classification Based on Decision Tree Algorithm for Machine Learning. *Journal of Applied Science and Technology Trends*, 2(01), 20–28. <https://doi.org/10.38094/jastt20165>
- Duarte, V., Zuniga-Jara, S., & Contreras, S. (2022). Machine Learning and Marketing: A Literature Review. *SSRN Electronic Journal*, 1–19. <https://doi.org/10.2139/ssrn.4006436>
- Fang, J. (2013). Why Logistic Regression Analyses Are More Reliable Than Multiple Regression Analyses. *Journal of Business and Economics*, 4(7), 2155–7950.
- Gkikas, D. C., Theodoridis, P. K., & Beligiannis, G. N. (2022). Enhanced Marketing Decision Making for Consumer Behaviour Classification Using Binary Decision Trees and a Genetic Algorithm Wrapper. *Informatics*, 9(2). <https://doi.org/10.3390/informatics9020045>
- Janzen, F. J., & Stern, H. S. (1998). Logistic regression for empirical studies of multivariate selection. *Evolution*, 52(6), 1564–1571. <https://doi.org/10.1111/j.1558-5646.1998.tb02237.x>
- Jiawei. (2012). *Data Mining* (Third Edit).
- Kusumaningtyas, R., Sholehah, I. M., & Kholifah, N. (2020). Peningkatan Kualitas

- Pembelajaran Guru Melalui Model dan Media Pembelajaran bagi Generasi Z. *Warta LPM*, 23(1), 54–62. <https://doi.org/10.23917/warta.v23i1.9106>
- Müller, A. C., & Guido, S. (2017). Introduction to Machine Learning with Python: a guide for data scientist. In *O'Reilly Media, Inc.*
- Santoso, G., & Triwijayati, A. (2018). Gaya Pengambilan Keputusan Pembelian Pakaian Secara Online pada Generasi Z Indonesia. *Jurnal Ilmu Keluarga Dan Konsumen*, 11(3), 231–242. <https://doi.org/10.24156/jikk.2018.11.3.231>
- Sperandeia, S., Leonardo, S. B., Marcelo, R. A., Reisd, A., & Basto, F. I. (2022). Assessing logistic regression applied to respondent-driven sampling studies: a simulation study with an application to empirical data. *INTERNATIONAL JOURNAL OF SOCIAL RESEARCH METHODOLOGY*, 8(1), 1–5.