

Perbandingan Algoritma Naive Dan Bayes Logistic Regression Untuk Penerimaan Siswa Baru (Studi Kasus Calon Siswa SMA Negeri 1 Brebes)

Ria Indah Fitria⁽¹⁾, Rizki Prasetyo Tulodo⁽²⁾, Nur Tulus Ujjianto⁽³⁾, Ali Sofian⁽⁴⁾
⁽¹⁾⁽²⁾⁽³⁾⁽⁴⁾Jurusan Informatika, Universitas Pancasakti Tegal
ria_indah@upstegal.ac.id⁽¹⁾

Abstrak

Di SMA (Sekolah Menengah Atas), siswa biasanya belajar dalam selama tiga tahun, meliputi kelas 10, 11, dan 12. Siswa biasanya berusia antara 15 hingga 18 tahun selama masa ini. SMA merupakan tahap penting dalam pendidikan di Indonesia karena siswa diharapkan memperoleh pengetahuan dan keterampilan yang lebih. SMA juga mempersiapkan siswa untuk ujian nasional yang biasanya diambil pada akhir tahun kelas 12. Hasil ujian nasional ini penting untuk menentukan kelayakan siswa untuk melanjutkan ke pendidikan tinggi, baik di universitas maupun institusi pendidikan vokasional. Pendaftaran di SMA memiliki beberapa kriteria, Misalnya 1) Apakah tempat tinggal calon siswa dengan sekolah jaraknya lebih dekat (zonasi), 2) nilai rata-rata dalam 3-5 semester pada raport, dan 3) Calon siswa memiliki prestasi baik itu non akademik atau akademik. Calon siswa yang memenuhi 3 kriteria diatas maka akan diterima di SMA. Dalam penelitian penulis membandingkan kedua algoritma dan metode adalah Algoritma *Naive Bayes* dan *logistic regresi*.

Kata Kunci Sekolah Menengah Atas (SMA), Pendaftaran, katagori, *Naive Bayes* dan *logistic regresi*

Pendahuluan

Sekolah Menengah Atas (SMA) merupakan Pendidikan yang khususnya mempersiapkan siswa lulusnya siap berkerja atau ingin melanjutkan kuliah. Kebanyakan siswa SMA akan siap berkerja saat lulus nanti atau melanjutkan ke jenjang kuliah.

SMA dirancang oleh pemerintah untuk mengembangkan keterampilan siswa, kemampuan dalam berkomunikasi siswa, sikap yang di tanamkan siswa agar saat berkerja memiliki sopan santun, disiplin siswa juga di latih agar berkerja lebih disiplin, dan apresiasi yang diperlukan dalam dunia kerja agar siswa mendapatkan pekerjaan yang layak sesuai dengan bidang mereka masing – masing.

Dalam dunia Pendidikan peran siswa disekolah sangatlah penting dalam memajukan sekolah. Jika kualitas muridnya menurun maka kualitas siswa yang diperoleh sekolah akan sangat menurun, oleh karena itu dalam penerimaan sekolah harus ada kriteria yang membantu dalam pemilihan. Kriteria yang ada dibuat oleh sekolah ada 3 kreteria, misalnya 1) Tempat jarak tinggal calon siswa dekat dengan sekolah, 2) nilai rata-rata raport dan 3) Prestasi non akademik atau akademik. Ketiga kriteria tersebut mambantu sekolah dalam menentukan siswa yang berkualitas.

Di sekolah bukan hanya nilai yang baik tetapi keterampilan juga sangan penting untuk terbentuknya siswa yang berkualitas. Maka dari itu sekolah memberikan fasilitas yang baik untuk menunjang system pembelajaran. Contohnya untuk jurusan IPA (Ilmu Pengetahuan Alam) sekolah memberikan laboratorium kimia yang berisi bahan atau praktek kimia untuk membantu siswa dalam pembuatan cairan atau larutan dalam pembelajaran.

Selain kualitas siswa yang baik. Sekolah juga meningkatkan mutu tenaga pendidik lebih baik. Banyak lulusan dari universitas terbaik yang berkerja di sekolah. Misalnya saja guru Bahasa Indonesia kebanyakan adalah lulusan dari sekolah negeri di Indonesia atau guru matematika yang memiliki banyak sekali piagam dalam mengajar sesuai dengan bidangnya dan masih banyak guru yang memiliki kualitas terbaik di sekolah SMA tersebut.

Untuk meningkatkan kualitas guru, sekolah mencari tenaga pengajar yang lulus Pendidikan S1 sampai dengan S3 untuk tenaga mengajar. Sekolah menginginkan siswa yang ada di sekolah memiliki kualitas yang baik agar nama sekolah akan menjadi baik dalam masyarakat. Selain siswa yang berkualitas tentu saja tenaga pengajar juga merupakan tenaga pengajar yang professional sesuai dengan bidangnya.

Landasan Teori

Pada Tahun 2016 Kurniawan[4], melakukan penelitian pada salah satu universitas di Indonesia salah satunya adalah di Fakultas Ilmu Komputer (FILKOM). Penelitian yang digunakan adalah hadoop dengan menerapkan naïve bayes yang dianggap mampu menghasilkan klasifikasi yang akurat, sehingga dapat mempermudah seorang dosen dalam pemilihan asisten praktikum dengan kualitas yang baik.

Pada Tahun 2016 Triowali Rosandy[5], Melakukan Penelitian salah satu algoritma klasifikasi yang digunakan adalah menggunakan metode Algoritma Naïve Bayes dan decision tree (c4.5). Tetapi kelemahan yang dihadapi pada kedua algoritma tersebut adalah lamanya waktu dan tingkat akurasi prediksi yang digunakan untuk melakukan prediksi. Hasil penelitian yang dihasilkan dari kedua algoritma adalah menganalisis prediksi pembiayaan yang di gunakan pada Bank di Indonesia.

Pada Tahun 2016 Ibnu[10] melakukan penelitian menggunakan metode regresi logistic untuk menghasilkan taksiran probabilitas munculnya peristiwa skor 1 pada variable kriteria. Regresi logistik memiliki model serupa dengan model regresi linier, yakni $p(Y=1)$ tetapi dengan konsep dan pemaknaan yang berbeda, yakni $\text{logit}(py) = \beta_0 + \beta X$. Untuk mengevaluasi apakah model regresi logistik cocok dengan data dapat dilakukan dengan Uji Wald, yang pada dasarnya merupakan χ^2 yang menguji apakah pengaruh masing-masing prediktor signifikan dalam memprediksi probabilitas terjadinya peristiwa pada variabel output.

Pada Tahun 2018 Wahyuni dkk [11], Melakukan Penelitian Peningkatan Akademik (PPA) merupakan salah satu beasiswa dengan jumlah peminat yang cukup banyak. Data dari bagian Akademik FMIPA UNNES menunjukkan bahwa jumlah peminat atau penyusul beasiswa PPA tahun 2015 di FMIPA UNNES sebesar 670 mahasiswa, Analisis data dalam penelitian ini menggunakan metode regresi logistic.

Metode Penelitian

Dalam Penelitian menggunakan jenis Perbandingan Kedua algoritma yang kemudian melakukan pengujian tingkat akurasi perbandingan algoritma untuk seleksi penerimaan calon siswa di SMA. Data tersebut diambil dari data calon siswa di SMA Negeri 1 Brebes Brebes.

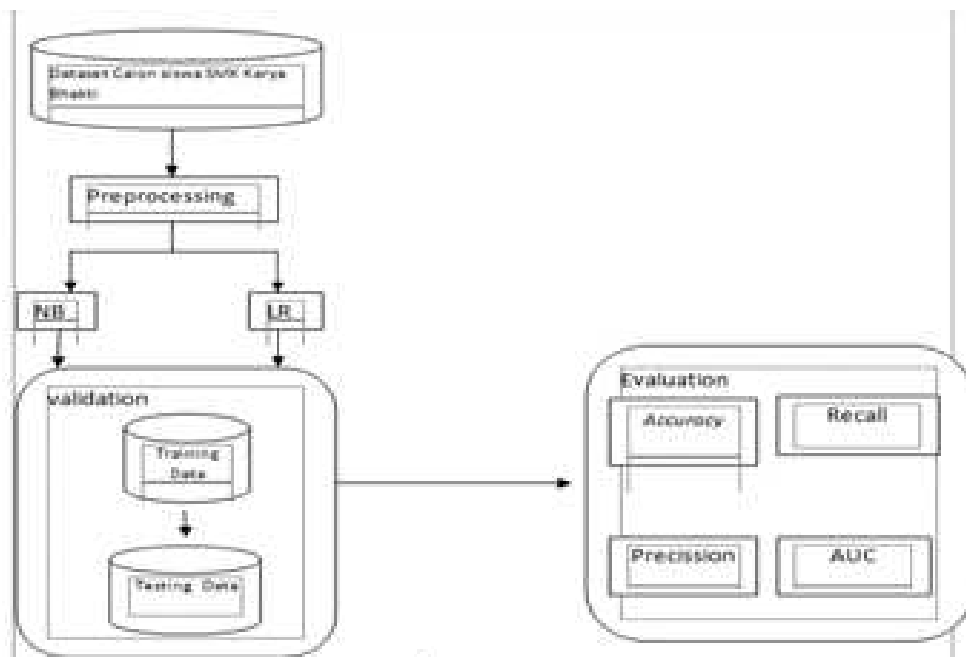
Penelitian ini menggunakan data dari SMA Negeri 1 Brebes. Data yang dibutuhkan adalah : jumlah data 36 siswa dan atribut yang dibutuhkan seperti nomer nis, L/P, jumlah un, rata-rata un, asal sekolah, alamat rumah, prestasi akademik. Atribut atribut tersebut diambil pada data calon penerimaan siswa di SMA. Data Kemudian dikelola untuk mendapatkan atribut mana saja yang relevan dan sesuai dengan format algoritma soft computing perbandingan sesuai dengan menggunakan tools rapid miner.

Sampel data untuk pengolahan pada rapid miner adalah data calon penerimaan siswa baru di SMA. SMA Negeri 1 Brebes pada tahun ajaran 2021/2022. Sampel data berjumlah 36 siswa terdiri dari 16 perempuan dan 20 laki – laki.

1. Metode yang di usulkan

Metode yang di usulkan adalah pembandingan antara logistic regression dan naïve bayes untuk mendeteksi permasalahan menyeleksi penerimaan siswa baru pada SMA Negeri 1 Brebes. Adapun kriteria adalah 1) Apakah tempat tinggal calon siswa dengan sekolah jaraknya lebih dekat (zonasi) , 2) nilai rata- rata dalam 3-5 semester pada raport , dan 3) Calon siswa memiliki prestasi

baik itu non akademik atau akademik Dari kedua algoritma tersebut yang memiliki akurasi yang lebih tinggi, maka akurasi yang lebih tinggi itu menunjukkan bahwa metodenya lebih baik.



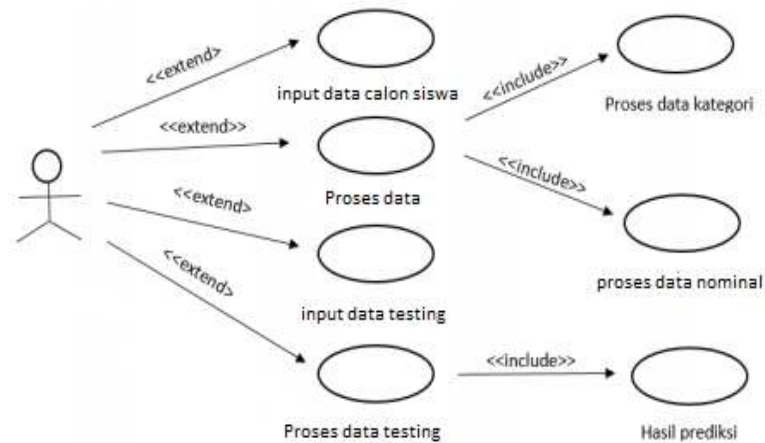
Gambar 1.1 Metode yang diusulkan

Dari gambar diatas menunjukkan bahwa metode yang diusulkan menunjukkan alur proses yang terjadi pada rapid miner. Dataset calon siswa SMA di prose pada rapid miner yang kemudian rapid miner memasukan algoritma yang akan di hitung. Disini menggunakan algoritma naïve bayes dan logistic regresi. Kedua algoritma tersebut divalidasi, validasi ada dua yang pertama training data gunanya untuk mengecek apakah data ada yang kosong atau terdapat missing dalam data dan kedua akan dilakukan testing data. Pada testing data disini dataset akan dihitung berdasarkan algoritma yang tadi di dibandingkan. Setelah di hitung maka akan muncul hasil evaluasi atau hasil akhir dari perhitungan. Evaluasi terdiri dari 4 yang antara lain: acsurasi, recall, precission, dan AUC.

Metode yang digunakan adalah untuk desain komparasi dua algoritma, dimana objek penelitian yang diambil dari calon siswa SMA Negeri 1 Brebes pada tahun 2021 dan tahun 2022 yang berupa data dari excel selanjutnya ditransportasikan ke aplikasi rapid miner. Sampel terdapat 36 siswa diantaranya perempuan 16 dan laki – laki 20 dengan menggunakan aplikasi rapidminer data diolah menggunakan dua metode algoritma.

Dalam pengujian metode komparasi algoritma dari hasil ekspresi. Menggunakan aplikasi rapidminer untuk dapat yang lebih baik. Dari hasil pengujian kedua algoritma akan diketahui algoritma manakah yang memiliki tingkat akurasi yang tinggi.

Evaluasi melakukan pengamatan dan sekaligus menganalisi hasil perbandingan algoritma. Validasi mengukur hasil deteksi klasifikasi permasalahan seleksi calon siswa menggunakan precision dan recall menakah hasilnya yang memiliki keakuratan atau presentasi yang tinggi.



Gambar 1.2 use case Naïve Bayes dan logistic regression

Gambar 1.2 adalah alur yang terjadi dalam aplikasi rapid miner untuk menjelaskan bagaimana algoritma naïve bayes dan logistic regression diproses kemudian database di validasi dan hasil dari validasi adalah nilai keakuratan dari metode tersebut. Rapid miner menganalisis hasil secara otomatis dan menghitung berapa presentasi yang di peroleh.

Hasil Penelitian Dan Pembahasan Kesimpulan

Dari penerapan menggunakan komparasi algoritma untuk klasifikasi seleksi penerimaan siswa di SMA dengan menggunakan rapid miner maka akan dapat dilihat bahwa algoritma Naïve Bayes menunjukkan hasil akurasi yang lebih baik dibandingkan dengan algoritma logistic regression

NO	NIS	L/P	jml un	rata 2 un	Asal Sekolah	Alamat Rumah	Prestasi akademik / non akademik	Status
1	13978	1	297	74.25	3	4	1	TERIMA
2	13977	1	266	66.5	3	2	2	TERIMA
3	13979	2	306	76.5	2	2	1	TERIMA
4	13980	1	306	76.5	4	3	2	TOLAK
5	13981	2	314	78.5	3	4	1	TERIMA
6	13982	1	320	80	2	2	1	TERIMA
7	13983	2	244	61.4	4	3	2	TOLAK
8	13984	2	297	74.25	4	4	1	TERIMA
9	13985	1	301	75.25	5	4	2	TOLAK
10	13986	2	281	70.25	2	2	1	TERIMA
11	13987	2	268	67	2	2	2	TOLAK
12	13988	2	345	86.25	3	4	2	TOLAK
13	13989	1	281	70.25	3	4	1	TERIMA
14	13990	1	270	67.5	4	5	1	TERIMA
15	13991	1	279	69.75	3	5	1	TERIMA
16	13992	1	300	75	3	5	1	TERIMA
17	13993	1	299	74.75	2	2	1	TERIMA
18	13994	1	314	78.5	5	1	2	TOLAK
19	13995	1	306	76.5	5	4	2	TOLAK
20	13996	2	294	73.5	2	4	1	TERIMA
21	13997	1	301	75.25	3	3	1	TERIMA
22	13998	1	305	76.25	2	1	1	TERIMA
23	13999	1	263	65.75	3	4	2	TOLAK
24	14000	1	292	73	2	3	1	TERIMA
25	14001	1	335	83.75	4	3	2	TOLAK
26	14002	1	252	63	3	5	1	TERIMA
27	14004	1	295	73.75	4	5	1	TERIMA

Gambar 3.1 Tabel dataset calon siswa SMA

Gambar di samping merupakan database calon siswa di SMA, data calon siswa ini nantinya akan dikelola dalam perhitungan algoritma *Naive Bayes*. Perhitungan yang nantinya hasil akurasi dari algoritma *Naive Bayes* akan dibandingkan dengan algoritma *Regresi Logistik*. Diantara dua algoritma tersebut manakah yang memiliki akurasi yang lebih besar.

A. Algoritma *Naive Bayes*

Algoritma *Naive Bayes* adalah salah satu metode klasifikasi yang populer dalam ilmu data dan pembelajaran mesin. Algoritma ini didasarkan pada teorema Bayes dengan asumsi bahwa fitur-fitur yang digunakan dalam klasifikasi adalah independen satu sama lain. Meskipun asumsi ini sering kali tidak terpenuhi di dunia nyata, *Naive Bayes* masih sering digunakan karena sederhana dan cepat serta memberikan hasil yang cukup baik dalam banyak kasus.

Berikut adalah langkah-langkah umum untuk mengimplementasikan algoritma *Naive Bayes*:

1. Pemilihan Model
2. Pemilihan Fitur
3. Pelatihan Model
4. Prediksi
5. Evaluasi Model [12]

accuracy: 97.50% +/- 7.50% (mikro: 97.14%)

	true TERIMA	true TOLAK	class precision
pred. TERIMA	23	0	100.00%
pred. TOLAK	1	11	91.67%
class recall	95.83%	100.00%	

Gambar 3.2 Tabel Accurassi Naïve Bayes

Gambar disamping adalah gambar akurasi dari algoritma naïve bayes. Pada gambar jelas tertera bahwa accurassi naïve bayes memiliki presentasi sebesar 97.50%.

Dalam pengukuran hasil penelitian menggunakan metode precision, recall, dan akurasi sedangkan untuk klasifikasi yang digunakan untuk melakukan pengukuran didapatkan dari hasil data calon siswa dari SMA sebanyak 36 siswa

Analisis Receiver Operating Characteristic (ROC) adalah metode yang umum digunakan untuk mengevaluasi performa model klasifikasi. ROC biasanya digunakan untuk mengukur seberapa baik model dapat membedakan antara kelas positif dan negatif dengan memvariasikan ambang batas (threshold) pengklasifikasian.

Tabel 3.1 Klasifikasi Pengujian Diagnostik

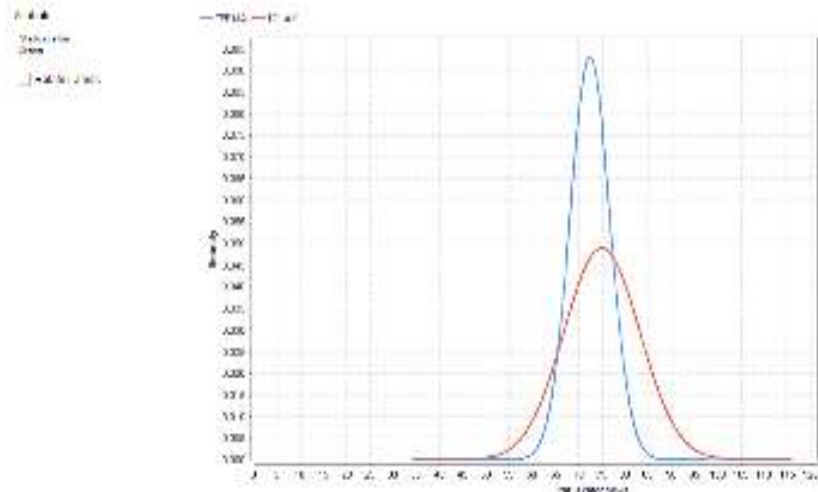
Nilai AUC	Klasifikasi
0.9 – 1.0	<i>Excellent Classification</i>
0.8 – 0.9	<i>Good Classification</i>
0.7 – 0,8	<i>Fair Classification</i>
0.6 – 0.7	<i>Poor Classification</i>
0.5 – 0.6	<i>Failure Classification</i>

Tabel 3.1 menunjukkan klasifikasi hasil perhitungan nilai AUC. Area Under the Curve (AUC) adalah metrik evaluasi yang umum digunakan dalam analisis kurva ROC (Receiver Operating Characteristic) untuk mengukur kinerja model klasifikasi. Nilai AUC berkisar antara 0 hingga 1, di mana semakin dekat ke 1 menunjukkan kinerja model yang lebih baik dalam membedakan antara kelas positif dan negatif.

```
PerformanceVector
PerformanceVector:
accuracy: 97.50% +/- 7.50% (mikro: 57.14%)
ConfusionMatrix:
True: TERIMA TOLAK
TERIMA: 23 0
TOLAK: 1 11
precision: 91.67% (positive class: TOLAK)
ConfusionMatrix:
True: TERIMA TOLAK
TERIMA: 23 0
TOLAK: 1 11
recall: 100.00% (positive class: TOLAK)
ConfusionMatrix:
True: TERIMA TOLAK
TERIMA: 23 0
TOLAK: 1 11
AUC (optimistic): 0.950 +/- 0.150 (mikro: 0.950) (positive class: TOLAK)
AUC: 0.900 +/- 0.200 (mikro: 0.900) (positive class: TOLAK)
AUC (pessimistic): 0.917 +/- 0.171 (mikro: 0.917) (positive class: TOLAK)
```

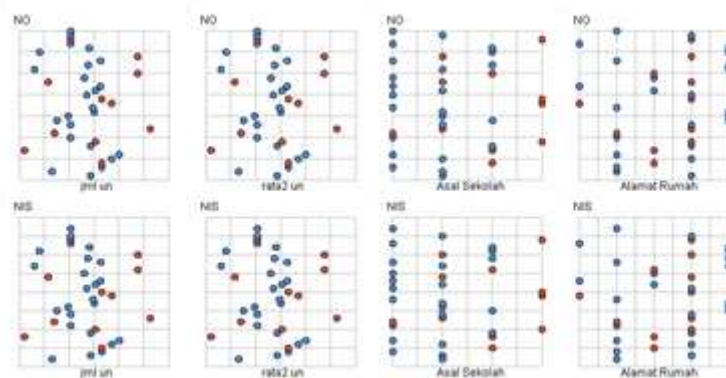
Gambar 3.3 Performance Vector Naïve Bayes

Gambar disamping adalah gambar hasil perhitungan dari rapat miner berdasarkan performance vectornya. Digambar tersebut di jelaskan berapa hasil recall, percession dan akursi dari algoritma Naïve bayes.



Gambar 3.4 Grafik presentasi berdasarkan status calon siswa

Gambar 3.4 merupakan gambar presentasi yang di tunjukan perhitungan naïve bayes. Dari gambar disamping bisa dilihat bahwa garis berwarna biru lebih dominan disbanding garis warna merah. Garis biru merupakan garis yang menandakan bahwa siswa tersebut di terima di SMA, sedangkan untuk garis merah adalah garis di tolak yang menandakan bahwa siswa tersebut di tolak / tidak di terima sebagai siswa di SMA



Gambar 3.5 Kluster diterima dan di tolak semua atribut pada algoritma Naïve Bayes

Gambar 3.5 adalah gambar kluster atau penggolongan semua atribut yang menyatakan di tolak dan diterima. Untuk lingkaran biru adalah diterima atau lolos menjadi siswa SMA, sedangkan untuk lingkaran merah adalah di tolak yang artinya tidak di terima sebagai siswa SMA. Dari kluster tersebut terlihat jelas bahwa tingkat diterima sangat besar dibandingkan dengan di tolak. Yang artinya kualitas dari calon siswa yang mendaftar merupakan siswa yang berprestasi dan memiliki jarak yang tidak begitu jauh dari tempat tinggalnya.

B. Algoritma Regresi Logistik

Algoritma Regresi Logistik adalah salah satu metode statistik yang digunakan untuk memodelkan hubungan antara satu atau lebih variabel independen (prediktor) dengan variabel dependen biner (target) untuk tujuan klasifikasi. Regresi logistik menghitung probabilitas bahwa suatu observasi akan termasuk dalam salah satu dari dua kategori atau kelas.

Berikut adalah langkah-langkah umum dalam menggunakan algoritma Regresi Logistik:

1. Pemilihan Variabel: Tentukan variabel independen yang akan digunakan untuk memprediksi variabel dependen biner. Pastikan variabel independen yang dipilih relevan dengan tujuan analisis.
2. Persiapan Data: Siapkan data Anda dengan melakukan pembersihan data, pengkodean variabel kategori menjadi variabel dummy jika diperlukan, dan pemisahan data menjadi set pelatihan (training set) dan set pengujian (test set).
3. Spesifikasi Model: Tentukan model regresi logistik yang sesuai dengan masalah Anda, termasuk apakah akan memasukkan interaksi antar variabel dan bagaimana menangani variabel yang tidak relevan atau kolinearitas.
4. Estimasi Parameter: Gunakan metode seperti Metode Maksimum Likelihood untuk mengestimasi parameter model regresi logistik dari data pelatihan.
5. Evaluasi Model: Evaluasi kinerja model Anda menggunakan metrik seperti akurasi, presisi, recall, F1-score, dan area di bawah kurva ROC (AUC). Gunakan set pengujian untuk menguji bagaimana model Anda berkinerja pada data yang belum pernah dilihat sebelumnya.
6. Penyetelan Model: Jika diperlukan, lakukan penyetelan model dengan mengubah parameter-model atau mencoba berbagai kombinasi variabel untuk meningkatkan kinerja model.
7. Penggunaan Model: Setelah Anda puas dengan kinerja model Anda, Anda dapat menggunakannya untuk melakukan prediksi pada data baru dengan menggunakan variabel independen yang relevan.[13].

Kelas	Jumlah	Prediksi	Akurasi
Kelas 0	23	23	100%
Kelas 1	1	1	100%
Kelas 2	0	0	0%

Gambar 3.6 Akurasi algoritma regresi logistik

Gambar disamping adalah gambar akurasi dari algoritma regresi logistik. Pada gambar jelas tertera bahwa akurasi naïve bayes memiliki presentasi sebesar 88.33%. Perhitungan akurasi algoritma regresi logistik adalah :

$$TP = 23$$

$$Fp = 3$$

$$Fn = 1$$

$$Tn = 8$$

$$\begin{aligned} \text{Precision} &= TP / (TP + FP) \\ &= 23 / (23 + 3) \\ &= 88.89 \end{aligned}$$

$$\begin{aligned} \text{Recall} &= TP / (TP + FN) \\ &= 23 / (23 + 1) \\ &= 72.73 \end{aligned}$$

$$\begin{aligned} \text{Akurasi} &= (TP + TN) / (TP + TN + FP + FN) \\ &= 0,8833 \end{aligned}$$

$$\begin{aligned} \text{Error rate} &= (FP + FN) / (TP + TN + FP + FN) \\ &= 0,021 \end{aligned}$$

Kelas	Actual	Predicted	Accuracy
Kelas 0	23	23	100%
Kelas 1	1	1	100%
Kelas 2	0	0	0%

Gambar disamping adalah gambar hasil perhitungan dari rapat miner berdasarkan performance vectornya. Digambar tersebut di jelaskan berapa hasil recall, percession dan akursi dari algoritma regresi logistik.

Kernel Model

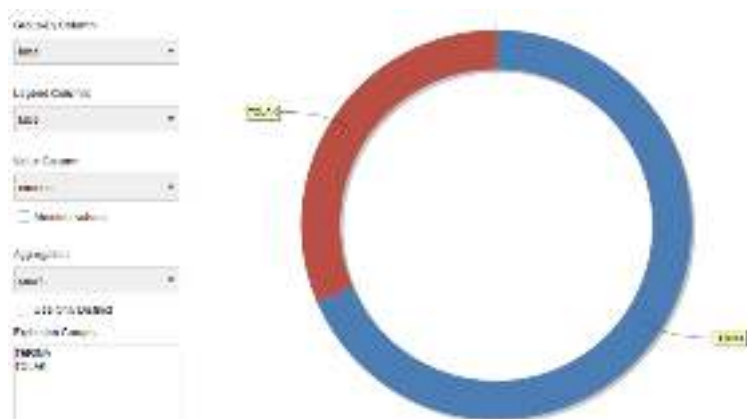
```
Total number of Support Vectors: 35  
Bias (offset): -1.506  
  
w[NO] = 0.122  
w[NIS] = 0.168  
w[L/P] = 0.327  
w[jml un] = 0.176  
w[rata2 un] = 0.176  
w[Asal Sekolah] = 0.593  
w[Alamat Rumah] = -0.034  
w[Prestasi akademik / non akademik] = 2.122
```

Gambar 3.8 Kernel Model

Gambar disamping adalah gambar hasil kernel model pada algoritma regresi logistik. Digambar tersebut di jelaskan model kernel perhitungan peratribut dalam datasetnya.

Hasilnya perhitungannya adalah sebagai berikut :

- NO = 0.122
- NIS = 0.168
- L/P = 0.327
- jml un = 0.176
- rata2 un = 0.176
- Asal Sekolah = 0.593
- Alamat Rumah = -0.034
- Prestasi akademik / non akademik = 2.122



Gambar 3.9 Diagram Ring presentasi berdasarkan status calon siswa

Gambar 3.9 merupakan gambar presentasi yang di tunjukan perhitungan regresi logistik. Dari gambar disamping bisa dilihat bahwa lingkaran berwarna biru lebih dominan disbanding lingkaran warna merah. lingkaran biru merupakan tanda bahwa siswa tersebut di terima di SMA,

- [7] A. N. Putri, "Penerapan Naive Bayesian Untuk Perankingan Kegiatan Di Fakultas Tik Universitas Semarang," *Simetris J. Tek. Mesin, Elektro dan Ilmu Komput.*, vol. 8, no. 2, p. 603, 2017.