

Data Mining Klasifikasi Untuk Memprediksi Status Keberlanjutan Polis Asuransi Kesehatan Dengan Algoritme Naïve Bayes

Jovansgha Avegad dan Arief Wibowo
Universitas Budi Luhur

Jl. Ciledug Raya, Petukangan Utara, Jakarta Selatan 12260
jovansghaavegad@gmail.com, arief.wibowo@budiluhur.ac.id

ABSTRAK

Di saat semakin kompetitifnya persaingan bisnis asuransi komersial kategori non-jaminan kesehatan nasional (JKN), membuat perusahaan asuransi umum dituntut memiliki inovasi dalam meningkatkan jumlah pelanggan atau nasabahnya. Di sisi lain, keputusan nasabah asuransi kesehatan terhadap keberlanjutan polis yang telah dimiliki apakah akan diperpanjang atau tidak diperpanjang, tidaklah mudah untuk diprediksi. Untuk dapat memprediksi keputusan nasabah asuransi kesehatan dalam status keberlanjutan polis maka diperlukan suatu metode analisis data pelanggan asuransi yang telah terdaftar. Penelitian ini menguraikan pemanfaatan algoritme Naïve Bayes sebagai salah satu algoritme terbaik dalam data mining klasifikasi untuk memprediksi keputusan nasabah asuransi terhadap polis yang dimiliki. Berdasarkan hasil pemodelan yang dilakukan, diketahui bahwa algoritme klasifikasi Naïve Bayes mampu memprediksi status keberlanjutan polis asuransi dengan akurasi mencapai 88,00%, presisi 89,19% dan recall 100%. Dengan hasil yang relatif baik ini maka dapat dilakukan upaya-upaya peningkatan pendapatan perusahaan asuransi kesehatan misalnya dengan menawarkan program promosi pembaharuan polis asuransi pada nasabah-nasabah yang diprediksi akan memperpanjang maupun yang tidak akan memperpanjang polis asuransi yang dimilikinya.

Kata Kunci : Data Mining, Klasifikasi, Naïve Bayes, Prediksi Pembaharuan Polis Asuransi Kesehatan

PENDAHULUAN

PT XYZ merupakan perusahaan yang bergerak dalam bidang asuransi dan memiliki nasabah yang cukup banyak. Data nasabah dari PT XYZ belum diolah sedemikian rupa untuk dimanfaatkan dalam prediksi pembaharuan polis asuransi.

Dengan proses Data Mining dari data nasabah, dapat ditemukan pola-pola ataupun hubungan keterkaitan tertentu antara data untuk menjadi informasi yang berharga [1]. Adapun tahapan mendapatkan informasi ini dapat diperoleh dalam tahap *Knowledge Discovery in Database (KDD)*.

Penelitian ini menguraikan bagaimana metode data mining algoritme Naïve Bayes atau yang disebut juga Bayesian Classification mampu digunakan untuk mengklasifikasi probabilitas status keberlanjutan polis asuransi kesehatan.

Naïve Bayes Classifier didasarkan pada teorema Bayes yang memiliki kemampuan klasifikasi serupa *decision tree* dan *neural network*. Selain itu, Naïve Bayes Classifier terbukti memiliki akurasi dan kecepatan yang tinggi saat diaplikasikan ke dalam database yang besar [2]. Adapun persamaan

atau formula yang digunakan dalam algoritme menurut teori Bayes tersebut adalah:

$$P(H | X) = (P(X | H)P(H)) / P(X) \quad (1)$$

Dimana :

X = Data dengan class yang belum diketahui

H = Hipotesis data X merupakan suatu class spesifik.

$P(H|X)$ = Probabilitas hipotesis H berdasarkan kondisi X (Posteriori Probability)

$P(H)$ = Probabilitas hipotesis H (Prior Probability)

$P(X|H)$ = Probabilitas X berdasarkan hipotesis H

$P(X)$ = Probabilitas dari X

Dalam proses evaluasi dari model yang dihasilkan, digunakan alat analisis *confusion matrix*, untuk mendapatkan informasi akurasi, presisi dan *recall* dari hasil pemodelan [3].

- a. Akurasi
Akurasi dalam klasifikasi adalah persentase ketepatan record data yang diklasifikasikan secara benar setelah dilakukan pengujian pada hasil klasifikasi.
- b. Presisi
Presisi didefinisikan sebagai rasio item relevan yang dipilih terhadap semua item yang terpilih. Presisi dapat diartikan sebagai kecocokan antara permintaan informasi dengan jawaban terhadap permintaan tersebut.
- c. Recall
Recall didefinisikan sebagai rasio dari item relevan yang dipilih terhadap total jumlah item relevan yang tersedia.
Hasil pengujian didapatkan dengan cara *holdout validation*, yaitu dengan membagi sebagian dataset menjadi data latih dan data uji.

METODE PENELITIAN

Tahapan metode penelitian yang digunakan yang direferensi berdasarkan tahapan dalam KDD adalah sebagai berikut:

- 1) *Data Selection*
Pemilihan (seleksi) data dari sekumpulan data operasional perlu dilakukan sebelum tahap penggalian informasi dalam KDD dimulai. Data hasil seleksi yang akan digunakan untuk proses data mining, disimpan dalam suatu berkas terpisah dari basis data operasional.
- 2) *Pre-Processing/Cleaning*
Proses *cleaning* antara lain membuang duplikasi data, memeriksa data yang tidak konsisten dan memperbaiki kesalahan pada data. Pada proses ini dilakukan juga proses enrichment, yaitu proses memperkaya data yang sudah ada dengan data atau informasi lain yang relevan dan diperlukan untuk KDD.
- 3) *Transformation*
Coding adalah proses transformasi pada data yang telah dipilih, sehingga data tersebut sesuai untuk proses data mining.
- 4) *Data Mining*
Data Mining adalah proses mencari pola atau informasi menarik dalam data

terpilih dengan menggunakan teknik atau metode tertentu.

5) *Interpretation/Evaluation*

Pola informasi yang dihasilkan dari proses data *mining* diterjemahkan menjadi bentuk yang lebih mudah dimengerti oleh pihak yang berkepentingan.

HASIL DAN PEMBAHASAN

Pada bagian ini akan dijelaskan seluruh tahapan yang telah dilakukan, mulai dari *data selection* hingga *interpretation* atau *evaluation*

Data Selection

Dataset yang diterima terdiri dari beberapa atribut yaitu:

- 1) No polis
- 2) Nama
- 3) Jenis kelamin
- 4) Usia
- 5) Pekerjaan
- 6) Tanggungan
- 7) Mulai asuransi
- 8) Plan
- 9) Premi tahunan
- 10) Periode asuransi
- 11) Riwayat penyakit
- 12) Rasio klaim
- 13) Perpanjangan

Dari atribut-atribut data yang dimiliki maka dilakukan pemilihan atribut sebagai acuan klasifikasi. Atribut yang dipilih sebagai berikut:

- 1) Jenis kelamin
- 2) Usia
- 3) Pekerjaan
- 4) Tanggungan
- 5) Plan
- 6) Premi tahunan
- 7) Riwayat penyakit
- 8) Rasio klaim
- 9) perpanjang

Pre-Processing/Cleaning

Tahap ini membersihkan data dari duplikasi dan data yang tidak konsisten, atau dikatakan tidak bersih, misalnya menghapus data yang tidak lengkap.

Transformation

Dilakukan transformasi pada beberapa data agar dapat diproses dengan metode

yang digunakan dalam data *mining*, transformasi data sebagai berikut:

Tabel 1. Transformasi data usia

Usia	Transformasi
1 – 28	KategoriUsia1
29 - 38	KategoriUsia2
39 - 48	KategoriUsia3
49 - 99	KategoriUsia4

Pembagian range usia ini diperoleh dari klasifikasi yang diberlakukan pada lokasi obyek penelitian. Transformasi lainnya yang dilakukan adalah:

Tabel 2. Transformasi data tanggungan

Tanggungan	Transformasi
0	Tanpa Tanggungan
1 - Seterusnya	Ada Tanggungan

Dengan jumlah tanggungan yang berisi nilai 0 dan seterusnya, maka dapat dilakukan kategorisasi tanggungan menjadi nilai ‘Ada tanggungan’ dan ‘Tanpa tanggungan’.

Premi tahunan juga ditransformasi berdasarkan median terhadap setiap plan yang ada. Transformasi juga dilakukan pada nilai premi sebagai berikut:

Tabel 3. Kelompok transformasi premi

Plan	Range (IDR)	
1	>822.600	<=822.600
2	>1.936.000	<=1.936.000
3	>3.335.400	<=3.335.400
4	>4.728.000	<=4.728.000
5	>4.728.000	<=4.728.000

Data Mining

Data mining adalah proses pencarian informasi didalam data dengan suatu metode yaitu Naïve Bayes.

Pemodelan ini menggunakan dataset faktual dari perusahaan asuransi kesehatan di DKI Jakarta pada tahun 2018 sejumlah 225 record yang dibagi menjadi 2/3 data latih dan 1/3 data uji sebagaimana metode validasi *holdout validation* [4].

Berikut ilustrasi penerapan analisis satu buah record dari dataset menggunakan algoritme Naïve Bayes:

- Jenis Kelamin : L
- Usia : dewasa awal
- Pekerjaan : Pegawai Swasta
- Tanggungan : ada
- Plan : D
- Premi tahunan : <=1936000
- Riwayat Penyakit : rendah
- Rasio Klaim : sedang
- Status : renewal

Berdasarkan persamaan dari penerapan Algoritme Naïve Bayes maka diperoleh hasil sebagai berikut:

- Kelas spesifik :
 - $P(\text{renewal}) = 0.91$
 - $P(\text{tidak renewal}) = 0.09$
- Jenis Kelamin :
 - $P(L|\text{renewal}) = 0.47$
 - $P(L|\text{tidak renewal}) = 0.38$
 - $P(L) = 0.47$
- Usia
 - $P(\text{dewasa awal}|\text{renewal}) = 0.26$
 - $P(\text{dewasa awal}|\text{tidak renewal}) = 0.46$
 - $P(\text{dewasa awal}) = 0.28$
- Pekerjaan
 - $P(\text{Pegawai Swasta}|\text{renewal}) = 0.58$
 - $P(\text{Pegawai Swasta}|\text{tidak renewal}) = 0.46$
 - $P(\text{Pegawai Swasta}) = 0.57$
- Tanggungan
 - $P(\text{ada}|\text{renewal}) = 0.41$
 - $P(\text{ada}|\text{tidak renewal}) = 0.15$
 - $P(\text{ada}) = 0.39$
- Plan
 - $P(D|\text{renewal}) = 0.55$
 - $P(D|\text{tidak renewal}) = 0.54$
 - $P(D) = 0.55$
- Premi Tahunan
 - $P(\leq 1936000|\text{renewal}) = 0.31$
 - $P(\leq 1936000|\text{tidak renewal}) = 0.08$
 - $P(\leq 1936000) = 0.29$
- Riwayat Penyakit
 - $P(\text{rendah}|\text{renewal}) = 0.42$
 - $P(\text{rendah}|\text{tidak renewal}) = 0.54$
 - $P(\text{rendah}) = 0.43$
- Rasio Klaim
 - $P(\text{sedang}|\text{renewal}) = 0.35$

- $P(\text{sedang}|\text{tidak renewal}) = 0.46$
- $P(\text{sedang}) = 0.36$

Probabilitas *renewal*:

$$\frac{P(X|\text{renewal})P(\text{renewal})}{P(X)} \quad (2)$$

Dimana X = Data dengan *class* yang belum diketahui.

$$P(X|\text{renewal}) = P(L|\text{renewal}) \times P(\text{dewasa awal}|\text{renewal}) \times P(\text{Pegawai Swasta}|\text{renewal}) \times P(\text{ada}|\text{renewal}) \times P(\text{D}|\text{renewal}) \times P(\leq 1936000|\text{renewal}) \times P(\text{rendah}|\text{renewal}) \times P(\text{sedang}|\text{renewal})$$

$$P(X|\text{renewal}) = 0.47 \times 0.26 \times 0.58 \times 0.41 \times 0.55 \times 0.31 \times 0.42 \times 0.35$$

$$P(X|\text{renewal}) = 0.000728324257$$

$$P(X) = P(L) \times P(\text{dewasa awal}) \times P(\text{Pegawai Swasta}) \times P(\text{ada}) \times P(\text{D}) \times P(\leq 1936000) \times P(\text{rendah}) \times P(\text{sedang})$$

$$P(X) = 0.47 \times 0.28 \times 0.57 \times 0.39 \times 0.55 \times 0.29 \times 0.43 \times 0.36$$

$$P(X) = 0.000223261186$$

$$P(\text{renewal}) = 0.91$$

$$(\text{renewal}|X) = \frac{0.000728324257 \times 0.91}{0.000722415602}$$

$$(\text{renewal}|X) = 0.917442913$$

Probabilitas tidak *renewal*:

$$\frac{P(X|\text{tidak renewal})P(\text{tidak renewal})}{P(X)} \quad (3)$$

Dimana X adalah data dengan *class* yang belum diketahui.

$$P(X|\text{tidak renewal}) = P(L|\text{tidak renewal}) \times P(\text{dewasa awal}|\text{tidak renewal}) \times P(\text{Pegawai Swasta}|\text{tidak renewal}) \times P(\text{ada}|\text{tidak renewal}) \times P(\text{D}|\text{tidak renewal}) \times P(\leq 1936000|\text{tidak renewal}) \times P(\text{rendah}|\text{tidak renewal}) \times P(\text{sedang}|\text{tidak renewal})$$

$$P(X|\text{tidak renewal}) = 0.38 \times 0.46 \times 0.46 \times 0.15 \times 0.54 \times 0.08 \times 0.54 \times 0.46$$

$$P(X|\text{tidak renewal}) = 0.00012942729$$

$$P(X) = P(L) \times P(\text{dewasa awal}) \times P(\text{Pegawai Swasta}) \times P(\text{ada}) \times P(\text{D}) \times P(\leq 1936000) \times P(\text{rendah}) \times P(\text{sedang})$$

$$P(X) = 0.47 \times 0.28 \times 0.57 \times 0.39 \times 0.17 \times 0.29 \times 0.43 \times 0.36$$

$$P(X) = 0.000223261186$$

$$P(\text{tidak renewal}) = 0.09$$

$$(\text{tidak renewal}|X) = \frac{0.00012942729 \times 0.09}{0.000722415602}$$

$(\text{tidak renewal}|X) = 0.0161243141$
Tahap akhir adalah membandingkan nilai probabilitas *renewal* dengan probabilitas tidak *renewal*:

$$(\text{renewal}|X) = 0.917442913$$

$$(\text{tidak renewal}|X) = 0.0161243141$$

$$(\text{renewal}|X) > (\text{tidak renewal}|X)$$

Dengan diketahuinya probabilitas *renewal* yang lebih besar, maka hasil analisis data dengan algoritme Naïve Bayes memberikan hasil prediksi bahwa keputusan nasabah terhadap polis asuransi yang dimiliki adalah akan diperbaharui (*renewal*).

Interpretation/Evaluation

Dilakukan proses evaluasi untuk merubah hasil informasi menjadi data yang lebih mudah dipahami, evaluasi dilakukan dengan metode *confusion matrix*.

Hasil pengujian terhadap 1/3 dataset sebagai data uji, mendapatkan nilai akurasi sebesar 88,0%, presisi 89,19% dan *recall* 100% sebagaimana tabel berikut:

Tabel 3. Hasil Confusion Matrix

Data Aktual	Hasil Prediksi	
	Renewal	Tidak Renewal
Renewal	66 (A)	1 (B)
Tidak Renewal	8 (C)	0 (D)

Perhitungan *Confusion Matrix*, sebagai berikut:

a. Perhitungan Akurasi

$$Akurasi = \frac{A + D}{A + B + C + D}$$

$$Akurasi = \frac{66 + 0}{66 + 8 + 1 + 0}$$

$$Akurasi = 0.88$$

b. Perhitungan Presisi

$$Presisi = \frac{A}{C + A}$$

$$Presisi = \frac{66}{8 + 66}$$

$$Presisi = 0.8919$$

c. Perhitungan Recall

$$Recall = \frac{A}{A + D}$$

$$Recall = \frac{66}{66 + 0}$$

$$Recall = 1$$

Dengan pengujian yang dilakukan didapatkan hasil akurasi yang relatif baik yaitu mencapai 88,0% disertai presisi sebesar 89,19% dan Recall sebesar 100%.

PENUTUP

Kesimpulan analisis dari pengujian prediksi perpanjangan polis asuransi dengan algoritme Naïve Bayes, memiliki akurasi 88,0%, presisi 89,19% dan *recall* 100%. Hasil yang relatif baik dan akurat ini memungkinkan bagi PT XYZ dapat melakukan upaya-upaya peningkatan pendapatan misalnya dengan menawarkan program perpanjangan polis asuransi pada nasabah-nasabah tertentu yang diprediksi akan memperpanjang polisnya, atau

memberikan penawaran menarik agar nasabah yang diprediksi tidak memperpanjang polisnya menjadi keputusan untuk memperbaharui/memperpanjang polis yang telah dimilikinya.

Saran

Saran untuk mengembangkan penelitian terhadap data nasabah asuransi lebih lanjut:

- Penelitian berikutnya dapat dikembangkan dengan menggunakan metode yang mungkin lebih baik, seperti C4.5 atau k-NN.
- Penelitian berikutnya dapat dikembangkan dengan metode validasi yang mungkin lebih baik, seperti *K-fold validation*.
- Pengumpulan data diharapkan lebih banyak dan bervariasi sehingga diharapkan akan meningkatkan nilai akurasi.

DAFTAR PUSTAKA

- N. Nuraeni, "Penentuan Kelayakan Kredit Dengan Algoritma Naïve Bayes Classifier: Studi Kasus Bank Mayapada Mitra Usaha Cabang PGC," *J. Tek. Komput. AMIK BSI*, vol. 3, no. 1, pp. 9–15, 2017.
- Y. V. Via, B. Nugroho, and A. Syafrizal, "Sistem Pendukung Keputusan Klasifikasi Tingkat Keganasan Kanker Payudara Dengan Metode Naive Bayes Classifier," *SCAN-Jurnal Teknol. Inf. dan Komun.*, vol. 10, no. 2, pp. 63–68, 2015.
- E. Mayadewi, P., & Rosely, "Prediksi Nilai Proyek Akhir Mahasiswa Menggunakan Algoritma Klasifikasi Data Mining," *Sist. Inf. Indones.*, no. November, pp. 2–3, 2015.
- T. Praningki and I. Budi, "Sistem Prediksi Penyakit Kanker Serviks Menggunakan CART, Naive Bayes, dan k-NN," *Creat. Inf. Technol. J.*, vol. 4, no. 2, p. 83, 2018.
- Bustami, "Penerapan Algoritma Naive Bayes Untuk Mengklasifikasi Data Nasabah Asuransi," *J. Inform.*, vol. 8, no. 1, pp. 884–898, 2014