

**PERBANDINGAN METODE REGRESI LOGISTIK BINER DAN  
*CLASSIFICATION AND REGRESSION TREES (CART)* UNTUK  
KLASIFIKASI DIAGNOSA PENYAKIT DIABETES MELLITUS (DM)**

**SKRIPSI**



**UIN SUNAN AMPEL  
S U R A B A Y A**

Disusun Oleh  
**PUJI WIDIARTI**  
**H72216062**

**PROGRAM STUDI MATEMATIKA  
FAKULTAS SAINS DAN TEKNOLOGI  
UNIVERSITAS ISLAM NEGERI SUNAN AMPEL  
SURABAYA**

**2020**

## PERNYATAAN KEASLIAN

Saya yang bertanda tangan di bawah ini,

Nama : PUJI WIDIARTI

NIM : H72216062

Program Studi : Matematika

Angkatan : 2016

Menyatakan bahwa saya tidak melakukan plagiat dalam penulisan skripsi saya yang berjudul " **PERBANDINGAN METODE REGRESI LOGISTIK BINER DAN CLASSIFICATION AND REGRESSION TREES (CART) UNTUK KLASIFIKASI DIAGNOSA PENYAKIT DIABETES MELLITUS (DM)** ".

Apabila suatu saat nanti terbukti saya melakukan tindakan plagiat, maka saya bersedia menerima sanksi yang telah ditetapkan.

Demikian pernyataan keaslian ini saya buat dengan sebenar-benarnya.

Surabaya, 23 Juli 2020

Yang menyatakan,



PUJI WIDIARTI  
NIM. H72216062

## LEMBAR PERSETUJUAN PEMBIMBING

Skripsi oleh

Nama : PUJI WIDIARTI

NIM : H72216062

Judul Skripsi : PERBANDINGAN METODE REGRESI LOGISTIK BINER  
DAN *CLASSIFICATION AND REGRESSION TREES (CART)*  
UNTUK KLASIFIKASI DIAGNOSA PENYAKIT  
DIABETES MELLITUS (DM)

telah diperiksa dan disetujui untuk diujikan.

Surabaya, 10 Juli 2020

Pembimbing



Wika Dianita Utami, M.Sc  
NIP. 19920610208012003

## PENGESAHAN TIM PENGUJI SKRIPSI

Skripsi oleh

Nama : PUJI WIDIARTI  
NIM : H72216062  
Judul Skripsi : PERBANDINGAN METODE REGRESI LOGISTIK BINER  
DAN *CLASSIFICATION AND REGRESSION TREES (CART)*  
UNTUK KLASIFIKASI DIAGNOSA PENYAKIT  
DIABETES MELLITUS (DM)

Telah dipertahankan di depan Tim Penguji  
pada tanggal 23 Juli 2020

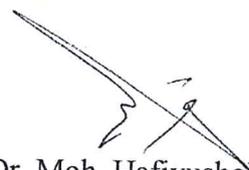
Mengesahkan,  
Tim Penguji

Penguji I



Wika Dianita Utami, M.Sc  
NIP. 199206102018012003

Penguji II



Dr. Moh. Hafiyusholeh, M.Si  
NIP. 198002042014031001

Penguji III



Putroue Keumala Intan, M.Si.  
NIP. 198805282018012001

Penguji IV



Nurissarah Ulinnuha, M.Kom  
NIP. 1999011022014032004

Mengetahui,

Plt. Dekan Fakultas Sains dan Teknologi  
UIN Sunan Ampel Surabaya



Dr. H. Evi Fatimatur Rusydiyah, M.Ag  
NIP. 197312272005012003



KEMENTERIAN AGAMA  
UNIVERSITAS ISLAM NEGERI SUNAN AMPEL SURABAYA  
PERPUSTAKAAN

Jl. Jend. A. Yani 117 Surabaya 60237 Telp. 031-8431972 Fax.031-8413300  
E-Mail: perpus@uinsby.ac.id

LEMBAR PERNYATAAN PERSETUJUAN PUBLIKASI  
KARYA ILMIAH UNTUK KEPENTINGAN AKADEMIS

Sebagai sivitas akademika UIN Sunan Ampel Surabaya, yang bertanda tangan di bawah ini, saya:

Nama : PUJI WIDIARTI  
NIM : H72216062  
Fakultas/Jurusan : SAINTEK/MATEMATIKA  
E-mail address : pujiwidiarti2@gmail.com

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Perpustakaan UIN Sunan Ampel Surabaya, Hak Bebas Royalti Non-Eksklusif atas karya ilmiah :  
 Skripsi  Tesis  Desertasi  Lain-lain (.....)  
yang berjudul :

PERBANDINGAN METODE REGRESI LOGISTIK BINER DAN CLASSIFICATION  
AND REGRESSION TREES (CART) UNTUK KLASIFIKASI DIAGNOSA PENYAKIT  
DIABETES MELLITUS (DM)

beserta perangkat yang diperlukan (bila ada). Dengan Hak Bebas Royalti Non-Eksklusif ini Perpustakaan UIN Sunan Ampel Surabaya berhak menyimpan, mengalih-media/format-kan, mengelolanya dalam bentuk pangkalan data (database), mendistribusikannya, dan menampilkan/mempublikasikannya di Internet atau media lain secara *fulltext* untuk kepentingan akademis tanpa perlu meminta ijin dari saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan atau penerbit yang bersangkutan.

Saya bersedia untuk menanggung secara pribadi, tanpa melibatkan pihak Perpustakaan UIN Sunan Ampel Surabaya, segala bentuk tuntutan hukum yang timbul atas pelanggaran Hak Cipta dalam karya ilmiah saya ini.

Demikian pernyataan ini yang saya buat dengan sebenarnya.

Surabaya, 23 Juli 2020

Penulis

(PUJI WIDIARTI)







3.3. Pengumpulan Data . . . . .	48
3.4. Tahap-tahap Penelitian . . . . .	50
3.5. <i>Flowchart</i> . . . . .	53
<b>IV HASIL DAN PEMBAHASAN . . . . .</b>	<b>55</b>
4.1. Deskriptif Data Diabetes Mellitus . . . . .	55
4.2. Analisis Klasifikasi Regresi Logistik Biner . . . . .	59
4.3. Analisis Klasifikasi Menggunakan <i>Classification And Regression Trees (CART)</i> . . . . .	65
4.4. Analisis Hasil . . . . .	99
<b>V PENUTUP . . . . .</b>	<b>105</b>
5.1. Simpulan . . . . .	105
5.2. Saran . . . . .	106
<b>DAFTAR PUSTAKA . . . . .</b>	<b>108</b>
<b>A Data Pasien Diabetes Mellitus di Rumah Sakit X . . . . .</b>	<b>112</b>
<b>B Data Training . . . . .</b>	<b>119</b>
<b>C Data Testing . . . . .</b>	<b>124</b>
<b>D Hasil Klasifikasi Penyakit Diabetes Mellitus Berdasarkan Data <i>Testing</i> Menggunakan Metode Regresi Logistik Biner . . . . .</b>	<b>127</b>
<b>E Pohon Klasifikasi Optimal CART . . . . .</b>	<b>128</b>
<b>F Hasil Klasifikasi Penyakit Diabetes Mellitus Berdasarkan Data <i>Testing</i> Menggunakan Metode CART . . . . .</b>	<b>129</b>
<b>G Skrip Program R . . . . .</b>	<b>130</b>









Terdapat beberapa faktor penyebab seseorang mengidap penyakit Diabetes Mellitus, antara lain faktor genetik, faktor lingkungan, faktor obesitas, faktor demografi dan lain sebagainya. Beberapa faktor risiko tersebut harus diperhatikan agar berdampak signifikan terhadap penurunan angka penderita Diabetes Mellitus serta penurunan kematian. Dengan demikian, perlu adanya program pengendalian Diabetes Mellitus, dengan cara mengendalikan faktor risiko tersebut (Misnadiarly, 2006).

Faktor risiko dapat dibedakan menjadi dua, yaitu faktor risiko dapat diubah dan tidak dapat diubah. Faktor risiko yang dapat diubah seperti pola hidup yang kurang baik misalnya kebiasaan merokok, kurang berolahraga dan mengkonsumsi makanan tidak sehat. Sementara itu, faktor risiko yang tidak dapat diubah misalnya faktor genetik, jenis kelamin dan usia.

Walaupun telah diketahui beberapa faktor yang berpengaruh terhadap pengidap Diabetes Mellitus, akan tetapi dalam menentukan seseorang terserang Diabetes Mellitus sangat sulit untuk ditentukan (WHO, 2006). Seiring dengan perkembangan teknologi dokter spesialis menggunakan rekam medis dan uji laboratorium untuk memperoleh data yang valid. Data yang didapatkan dari hasil rekam medis tersebut selanjutnya akan didiagnosa oleh dokter apakah pasien mengidap Diabetes Mellitus. Berdasarkan data tersebut dapat dibentuk model matematis, yaitu dengan metode klasifikasi (Widagdo, 2010).

Metode klasifikasi salah satunya dapat diselesaikan menggunakan metode statistik. Pada penyelesaian masalah klasifikasi (pengelompokkan) perlu diperhatikan dalam menentukan metode klasifikasi yang tepat, misalkan ingin mengklasifikasikan pasien yang mengidap penyakit Diabetes Mellitus dan tidak mengidap Diabetes Mellitus. Apabila mengklasifikasikan pasien pengidap penyakit Diabetes Mellitus ke dalam kelompok pasien yang tidak mengidap

penyakit Diabetes Mellitus maka akan menimbulkan suatu kesalahan yang fatal. Pengklasifikasian suatu objek dapat menggunakan pendekatan metode parametrik dan nonparametrik (Widagdo, 2010).

Salah satu metode statistik dengan pendekatan parametrik adalah Regresi Logistik Biner. Analisis Regresi merupakan salah satu metode statistik yang digunakan untuk mempelajari hubungan antara variabel independen dengan variabel dependen. Regresi Logistik merupakan salah satu metode yang digunakan untuk menganalisa data yang mempunyai variabel dependen berupa data kategorik (Margasari, 2014).

Regresi Logistik tidak mengasumsikan hubungan linier antar variabel independen dan dependen, dikarenakan bentuk variabel dependen berupa kategorik. Variabel dependen dengan dua kategorik (biner) yang bertolak belakang dapat dilakukan pengujian statistika menggunakan metode Regresi Logistik Biner. Variabel dependen biasanya diprediksi dengan nilai 0 dan 1, dimana 0 menyatakan bahwa pasien negatif Diabetes Mellitus dan 1 menyatakan pasien positif Diabetes Mellitus (Rumaendra, 2016).

Sementara itu, klasifikasi dapat menggunakan metode pendekatan nonparametrik. Pendekatan nonparametrik tidak bergantung terhadap asumsi tertentu, sehingga dapat memperoleh fleksibilitas yang lebih besar untuk menganalisa suatu data. Akan tetapi, tetap memiliki tingkat akurasi yang tinggi serta mudah dalam penerapannya. Salah satu metode statistik untuk klasifikasi dengan pendekatan nonparametrik adalah metode *Classification And Regression Trees* (CART). Pada tahun 1984, metode tersebut diperkenalkan oleh Leo Breiman, Richard A. Olshen, Jerome H. Friedman serta Charles J. Stone. Keempat ilmuwan tersebut memperkenalkan metode *Classification And Regression Trees* (CART) sebagai metode pohon klasifikasi dan regresi. CART dikembangkan untuk

menganalisa klasifikasi, baik variabel dependen kontinu maupun kategorik. CART akan menghasilkan *output* pohon klasifikasi (*classification trees*) apabila variabel Y berupa kategorik dan menghasilkan *output* pohon regresi (*regression trees*) apabila variabel Y berupa kontinu atau numerik (Breiman dkk., 1984).

Metode CART mempunyai kelebihan struktur data yang digunakan dapat dilihat secara visual, dapat mengeksplorasi struktur data yang kompleks dan bersifat nonparametrik sehingga tidak membutuhkan asumsi tertentu yang sering tidak terpenuhi oleh suatu data, serta proses pengklasifikasian lebih mudah dilakukan dengan cara menelusuri pohon klasifikasi yang dihasilkan (Breiman dkk., 1993).

Berdasarkan penelitian yang telah dilakukan oleh Wella Rumaenda (2016) dengan judul Perbandingan Klasifikasi Penyakit Hipertensi Menggunakan Regresi Logistik Biner dan Algoritma C4.5. Dalam penelitian tersebut bertujuan untuk membandingkan ketepatan klasifikasi antara metode Regresi Logistik Biner dengan Algoritma C4.5. Berdasarkan metode yang digunakan hasil klasifikasi menggunakan metode Regresi Logistik Biner diperoleh ketepatan klasifikasi sebesar 72,5352% sedangkan menggunakan Algoritma C4.5 didapatkan ketepatan klasifikasi sebesar 64,0845% (Rumaendra, 2016).

Penelitian lain yang telah dilakukan oleh R. Lestawati, Rais dan I. T. Utami (2016) dengan judul Perbandingan antara Metode CART (*Classification And Regression Tress*) dan Regresi Logistik (*Logistic Regression*) dalam Mengklasifikasikan Pasien Penderita DBD (Demam Berdarah *Dengue*). Penelitian tersebut bertujuan untuk memperoleh ketepatan hasil klasifikasi antara metode CART dan Regresi Logistik serta variabel yang signifikan terhadap penyakit DBD. Berdasarkan metode yang digunakan hasil klasifikasi menggunakan metode CART diperoleh ketepatan klasifikasi sebesar 76,3% dengan variabel yang signifikan

yaitu hepatomegali, epitaksis, melena dan diare sedangkan hasil klasifikasi menggunakan metode Regresi Logistik Biner diperoleh ketepatan klasifikasi sebesar 76,7% dan variabel yang signifikan yaitu hepatomegali (Lestawati dkk., 2018).

Penelitian lain yang telah dilakukan oleh Rifqy Marwah Akhsanti, Widyanti Rahayu dan Vera Maya Santi (2018) dengan judul Klasifikasi Diagnosis Penyakit Kanker Payudara dengan Pendekatan Regresi Logistik Biner dan Metode *Classification And Regression Tress (CART)*. Penelitian tersebut bertujuan untuk mengetahui faktor pengaruh timbulnya kanker payudara yang diklasifikasikan menurut dua kategorik, yaitu jinak dan ganas serta mengetahui ketepatan hasil klasifikasi dari dua metode tersebut. Berdasarkan penelitian, hasil yang diperoleh dengan analisis Regresi Logistik Biner faktor yang berpengaruh signifikan terhadap hasil diagnosis kanker payudara adalah usia dan riwayat keluarga penderita kanker (RKPK) serta ketepatan hasil klasifikasi yang diperoleh sebesar 90,5%. Sementara itu, metode CART menghasilkan pohon klasifikasi optimum dengan empat simpul terminal dan ketepatan hasil klasifikasi yang diperoleh sebesar 93% (Akhsanti dkk., 2018).

Berdasarkan latar belakang tersebut, tidak selamanya akurasi dari metode CART lebih baik daripada Regresi Logistik Biner untuk klasifikasi suatu kasus, begitupula sebaliknya. Oleh karena itu, dalam penelitian ini akan dilakukan klasifikasi penyakit Diabetes Mellitus menggunakan perbandingan metode Regresi Logistik Biner dan CART. Dengan demikian, dalam penelitian ini penulis mengambil judul “PERBANDINGAN METODE REGRESI LOGISTIK BINER DAN *CLASSIFICATION AND REGRESSION TREES (CART)* UNTUK KLASIFIKASI DIAGNOSA PENYAKIT *DIABETES MELLITUS (DM)*”.



















## 2.4.

**Tabel 2.4 Kriteria Indeks Massa Tubuh (IMT)**

<b>Kategori</b>	<b>Indeks Massa Tubuh (IMT)</b>
BB Normal	$\leq 33$
BB Tidak Normal	$\geq 34$

Berdasarkan Tabel 2.4, hasil pengukuran IMT yang masuk kategori tidak normal (obesitas) perlu untuk diwaspadai. Obesitas merupakan salah satu faktor risiko yang berperan penting terhadap penyakit Diabetes Mellitus. Seseorang dengan obesitas mempunyai masukan kalori yang berlebih. Sel beta pada kelenjar pankreas akan mengalami kelelahan sehingga tidak mampu untuk memproduksi insulin yang cukup untuk mengimbangi masukan kalori yang berlebih. Dengan demikian, kadar gula darah akan tinggi sehingga akan menjadi Diabetes Mellitus (Kaban & Sempakata, 2007).

## 5. Pola Makan

Seseorang dengan gaya hidup tidak baik, seperti pola makan tidak normal. Pola makan yang kurang baik misalnya mengkonsumsi alkohol, merokok dan konsumsi lemak berlebihan dapat memicu seseorang mengidap Diabetes Mellitus (Kemenkes, 2010).

Dengan pemanfaatan diagnosa serta data laboratorium, teknologi informasi dapat memberikan solusi untuk menyelesaikan permasalahan klasifikasi Diabetes Mellitus. Cabang ilmu data mining yang secara khusus berurusan mengenai informasi, penyimpanan, penarikan serta penggunaan data dapat memecahkan permasalahan dan juga pengambilan keputusan yang berasal dari dunia media (ADA, 2009).

### 2.3. Klasifikasi

Pada klasifikasi, terdapat target pada variabel dependen berupa kategorik. Sebagai contoh, pengklasifikasian diagnosa Diabetes Mellitus dapat dibagi dalam dua kategorik, yaitu positif Diabetes Mellitus dan negatif Diabetes Mellitus (Kusrini & Luthfi, 2009). Proses klasifikasi akan menggunakan data berupa kategorik yang bernilai diskrit/kontinu. Menurut Goronescu, dalam proses klasifikasi didasarkan pada empat komponen dasar yang sangat penting, komponen tersebut adalah sebagai berikut (Goronescu, 2011):

#### 1. *Class*

Variabel dependen dari model, merupakan variabel berupa kategorik yang merepresentasikan label terhadap objek setelah klasifikasi. Contohnya terdapat *class* pada klasifikasi penyakit Diabetes Mellitus, yaitu *class* positif Diabetes Mellitus dan *class* negatif Diabetes Mellitus (Goronescu, 2011).

#### 2. Prediktor

Variabel independen dari model, merupakan faktor-faktor yang mempengaruhi variabel dependen. Contohnya adalah data gejala dan hasil uji laboratorium mengenai faktor-faktor yang mempengaruhi Diabetes Mellitus, misalnya faktor genetik (Goronescu, 2011).

#### 3. *Training* dataset

Kumpulan data yang terdiri dari dua komponen di atas yang digunakan untuk melatih model dalam mengenali *class* yang sesuai berdasarkan variabel independen yang tersedia (Goronescu, 2011).

#### 4. *Testing* dataset

Kumpulan data baru yang akan diklasifikasikan menggunakan model yang

telah dibentuk sebelumnya. Dengan demikian, akan memperoleh hasil akurasi klasifikasi dan dapat dievaluasi (Goronescu, 2011).

#### 2.4. Regresi Logistik Biner

Pada dasarnya analisis regresi merupakan suatu ilmu mengenai hubungan antara variabel dependen dengan satu atau lebih variabel independen, dengan maksud untuk memprediksi dan mengestimasi nilai-nilai variabel dependen berdasarkan nilai variabel independen yang telah diketahui (Ghozali & Imam, 2005).

Regresi Logistik Biner merupakan salah satu metode analisis data yang digunakan untuk mengetahui keterkaitan antara beberapa variabel independen ( $x$ ) terhadap variabel dependen ( $y$ ) yang mempunyai dua kategori (biner) dimana variabel dependen ( $y$ ) bernilai salah dan benar (Hosmer & Lemeshow, 2000).

Metode yang cocok untuk mengklasifikasikan kejadian dengan variabel dependen ( $y$ ) yang mempunyai dua kategori yaitu menggunakan metode regresi logistik, misalnya dalam mengklasifikasikan pasien apakah menderita penyakit Diabetes Mellitus atau tidak dengan memperhatikan beberapa faktor yang mempengaruhinya. *Output* dari variabel dependen ( $y$ ) yang akan terbentuk dari pasien tersebut yaitu negatif Diabetes Mellitus atau positif Diabetes Mellitus (Hosmer & Lemeshow, 2000).

Apabila terdapat variabel independen sebanyak  $k$ , maka probabilitas untuk memperoleh "NEGATIF" ( $y = 0$ ) dapat dinyatakan dalam bentuk  $P(Y = 0|x) = \pi(x)$ . Sementara itu, probabilitas untuk memperoleh "POSITIF" ( $y = 1$ ) dapat dinyatakan dalam bentuk  $P(Y = 1|x) = \pi(x)$ . Variabel ( $x$ ) merupakan variabel independen yang dapat bersifat kualitatif, misalkan  $x = 0$  atau  $x = 1$  dan seterusnya. Bentuk umum fungsi regresi logistik dengan variabel independen



































disebut dengan istilah *binary recursive partitioning*. Proses *binary* yaitu setiap *parent node* akan selalu mengalami pemecahan tepat dua *child node*. Proses *recursive* berarti bahwa pada proses *binary* tersebut akan diulang kembali disetiap *child nodes* sebagai hasil pemecahan terdahulu, sehingga *child nodes* tersebut sekarang menjadi *parent nodes*. Proses pemecahan ini akan dilakukan perulangan sampai tidak dapat melakukan pemecahan berikutnya. Sementara itu, proses *partitioning* berarti bahwa *learning sample* yang dimiliki akan dipecah ke dalam partisi-partisi yang lebih kecil (Breiman dkk., 1984).

*Splitting criterion* (kriteria pemecahan) didasarkan pada nilai-nilai dari variabel independen yang digunakan. Misalkan variabel dependen  $Y$  dengan berskala kategorik dan variabel-variabel independen  $X_1, X_2, \dots, X_k$ . Proses *binary recursive partitioning* dapat diilustrasikan sebagai proses penyekatan atau pembagian dari ruang berdimensi  $k$  dari variabel-variabel  $X$  kedalam sekatan-sekatan yang berbentuk persegi panjang dan tidak saling bertumpang tindih. Pertama, pilih salah satu variabel independen  $X_k$  dan nilai  $X_k$  misalkan  $S_1$  terpilih untuk memecah ruang berdimensi  $k$  kedalam dua bagian. Bagian pertama adalah sekatan berisi objek-objek yang mana  $X_k \leq S_1$  (Breiman dkk., 1984).

Sementara itu, bagian lainnya adalah sekatan yang berisi objek-objek dengan nilai  $X_k > S_1$ . Selanjutnya, masing-masing dari sekatan tersebut dipecah kembali dengan cara yang sama dengan variabel independen (variabel independen yang terpilih dapat  $X_k$  kembali atau yang lainnya) dengan nilai tertentu. Proses ini akan terus berlanjut hingga diperoleh sekatan-sekatan yang lebih kecil dengan berisikan objek-objek yang homogen atau seragam. Homogen yang dimaksud adalah objek-objek yang terdapat dalam sekatan tersebut merupakan anggota satu *class* yang sama. Akan tetapi, pada kenyataannya keadaan seperti tersebut tidaklah mutlak selalu diperoleh. Proses *splitting* akan terus berlanjut hingga diperoleh

















Setelah dilakukan pemangkasan bagian pohon yang kurang penting, maka akan diperoleh pohon klasifikasi optimal. Metode pemangkasan pohon pada sebelumnya memperoleh urutan *subtree*  $T_1 > T_2 > \dots > \{t\}$ , karena pohon klasifikasi yang diperoleh berjumlah cukup banyak, maka permasalahan yang timbul adalah bagaimana cara untuk menentukan pohon klasifikasi yang optimal (Breiman dkk., 1993). Apabila menggunakan Persamaan (2.21)

$$R(T) = \sum_{t \in \bar{T}} r(t)P(t) = \sum_{t \in \bar{T}} R(t)$$

maka  $T_1$  akan dipilih sebagai pohon klasifikasi optimal, karena nilai *resubstitution estimate* dari  $T_1$  paling kecil. Oleh karena itu, metode *resubstitution estimate* merupakan metode yang bias untuk mengestimasi nilai *true misclassification cost*. Ada dua metode tak bias untuk mengestimasi nilai *true misclassification cost* yaitu penduga sampel uji *test sample estimate* dan penduga validasi silang lipat V (*cross validation V-fold estimate*) (Breiman dkk., 1993).

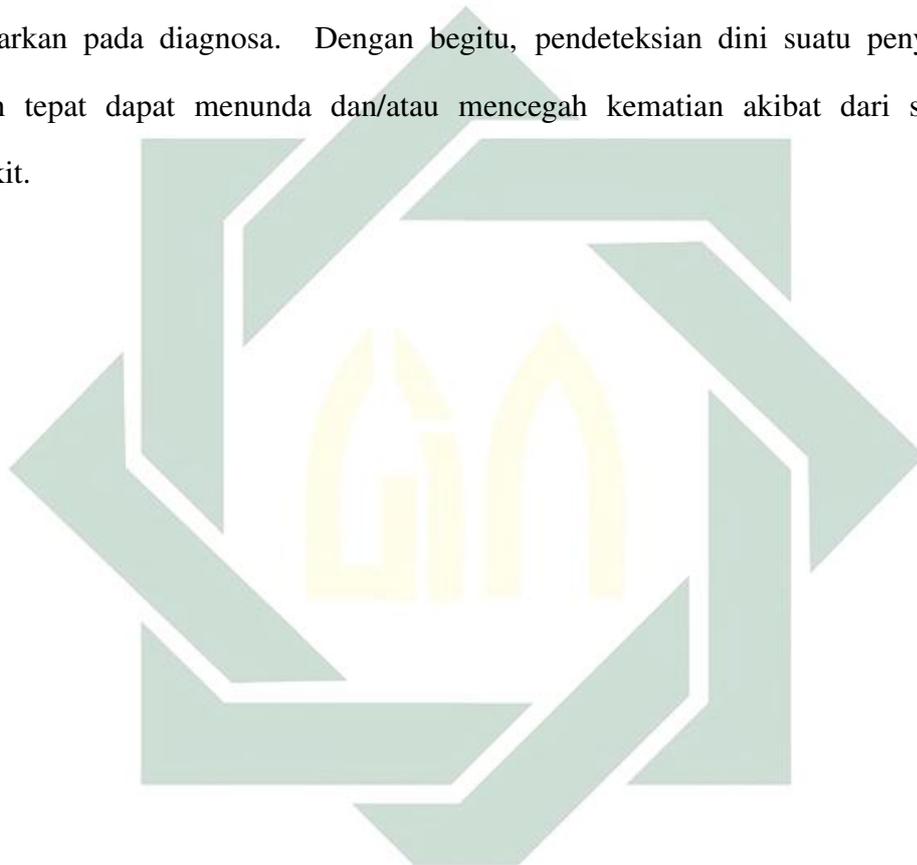
*Test sample estimate* digunakan apabila ukuran data dengan jumlah besar ( $L$ ). Kemudian data ( $L$ ) tersebut dibagi menjadi dua bagian  $L_1$  dan  $L_2$ . Misalkan  $L_1$  sebanyak  $N_1$  objek dan  $L_2$  sebanyak  $N_2$  objek. Selanjutnya, dibuat  $T_{max}$  menggunakan  $L_2$  dan dilakukan pemangkasan hingga diperoleh  $\{T_k\} = T_1 > T_2 > \dots > \{t\}$ .  $L_2$  digunakan untuk memprediksi pada masing-masing *tree* dan dihitung banyak objek yang terjadi kesalahan klasifikasi (*misclassification*) (Breiman dkk., 1993).

Sementara itu, *cross validation* digunakan jika ukuran data berjumlah kecil. Dalam *V-fold cross validation learning sample L* dibagi secara acak dalam V bagian terpisah, sebanyak  $V-1$  bagian sebagai *learning sample* dan sisa ( $V-(V-1)$ ) sebagai *test sample*. Setiap bagian terdapat objek dengan jumlah yang





walaupun bagaimanapun penyakit tersebut, asalkan tetap berusaha dan berdoa kepada-Nya. Hal ini dijelaskan juga pada H.R. Muslim yang artinya *"Setiap penyakit ada obatnya. Apabila ditemukan obat yang tepat untuk suatu penyakit, akan sembuhlah penyakit itu dengan izin Allah 'azza wajalla"*. Bentuk usaha yang dapat dilakukan manusia yaitu dengan melakukan perawatan medis yang tepat berdasarkan pada diagnosa. Dengan begitu, pendeteksian dini suatu penyakit dengan tepat dapat menunda dan/atau mencegah kematian akibat dari suatu penyakit.































Tabel 4.5 Nilai  $Z_{value}$  dan  $P_{value}$ 

Variabel	$Z_{value}$	$P_{value}$ / Sig.
(Intercept)	-5,264	1,41e-07 ***
Usia ( $X_1$ )	2,137	0,0326 *
Keturunan ( $X_2$ )	4,327	1,51e-05 ***
Gula Darah ( $X_3$ )	5,651	1,59e-08 ***
Obesitas ( $X_4$ )	4,117	3,83e-05 ***

Berdasarkan hasil pengujian secara individu, diperoleh nilai  $P_{value}$  masing-masing variabel independen  $X_1 = 0,0326$ ,  $X_2 = 1,51e - 05$ ,  $X_3 = 1,59e - 08$  dan  $X_4 = 3,83e - 05$ . Masing-masing nilai  $P_{value}$  tersebut  $< (\alpha = 0,05)$ . Oleh karena itu, secara individu variabel independen  $X_1$ ,  $X_2$ ,  $X_3$  dan  $X_4$  mempunyai pengaruh terhadap variabel diabetes ( $Y$ ).

### 3. Diagnosa Diabetes Mellitus

Setelah diperoleh suatu model Regresi Logistik Biner dan masing-masing uji terpenuhi, selanjutnya akan dilakukan diagnosa penyakit Diabetes Mellitus dari data *testing* yang telah disiapkan. Data tersebut diuji menggunakan model yang diperoleh sebelumnya.

Misalkan untuk mengetahui seseorang menderita penyakit DM dengan ciri-ciri muda, tidak mempunyai riwayat keturunan DM, kadar gula darah normal dan memiliki tingkat kegemukan (obesitas) normal. Maka dapat dihitung menggunakan model Regresi Logistik Biner yang telah dibentuk sebagai berikut:

$$\begin{aligned}\pi(x_i) &= \frac{e^{(-6,9683+1,6115X_1+5,6213X_2+3,9764X_3+4,6913X_4)}}{1+e^{(-6,9683+1,6115X_1+5,6213X_2+3,9764X_3+4,6913X_4)}} \\ &= \frac{e^{(-6,9683+1,6115(0)+5,6213(0)+3,9764(0)+4,6913(0))}}{1+e^{(-6,9683+1,6115(0)+5,6213(0)+3,9764(0)+4,6913(0))}}\end{aligned}$$

$$= \frac{e^{-6,9683}}{1+e^{-6,9683}}$$

$$= 0,000940367$$

Berdasarkan model tersebut seseorang dengan ciri-ciri muda, tidak mempunyai riwayat keturunan DM, kadar gula darah normal dan tingkat kegemukan (obesitas) normal sebesar  $0,000940367 \leq 0,5$ . Sehingga seseorang dengan ciri-ciri tersebut dinyatakan negatif Diabetes Mellitus.

Dengan demikian, akan diperoleh diagnosa pasien positif dan negatif Diabetes Mellitus pada Tabel 4.6 sebagai berikut:

No.	Aktual	Probabilitas	Prediksi	No.	Aktual	Probabilitas	Prediksi	No.	Aktual	Probabilitas	Prediksi
1	Negatif	0.09305	Negatif	25	Positif	0.84546	Positif	49	Positif	0.84546	Positif
2	Negatif	0.20636	Negatif	26	Positif	0.99934	Positif	50	Positif	0.09305	Negatif
3	Negatif	0.20094	Negatif	27	Positif	0.84546	Positif	51	Positif	0.96480	Positif
4	Negatif	0.33951	Negatif	28	Positif	0.99934	Positif	52	Positif	0.84546	Positif
5	Negatif	0.00094	Negatif	29	Positif	0.84546	Positif	53	Negatif	0.09305	Negatif
6	Negatif	0.33951	Negatif	30	Positif	0.99987	Positif	54	Negatif	0.33951	Negatif
7	Negatif	0.00094	Negatif	31	Positif	0.99301	Positif	55	Negatif	0.00469	Negatif
8	Negatif	0.09305	Negatif	32	Positif	0.99301	Positif	56	Negatif	0.20636	Negatif
9	Negatif	0.96592	Positif	33	Positif	0.96480	Positif	57	Negatif	0.20636	Negatif
10	Negatif	0.04779	Negatif	34	Positif	0.84546	Positif	58	Negatif	0.84546	Positif
11	Negatif	0.04779	Negatif	35	Positif	0.99934	Positif	59	Negatif	0.20636	Negatif
12	Negatif	0.04779	Negatif	36	Positif	0.99934	Positif	60	Positif	0.84546	Positif
13	Negatif	0.84546	Positif	37	Positif	0.96480	Positif	61	Positif	0.98581	Positif
14	Negatif	0.20094	Negatif	38	Positif	0.99934	Positif	62	Positif	0.96592	Positif
15	Negatif	0.84546	Positif	39	Positif	0.96480	Positif	63	Positif	0.96480	Positif
16	Positif	0.99934	Positif	40	Positif	0.96592	Positif	64	Positif	0.99934	Positif
17	Positif	0.96592	Positif	41	Positif	0.84546	Positif	65	Positif	0.99987	Positif
18	Positif	0.99934	Positif	42	Positif	0.99301	Positif	66	Positif	0.84546	Positif
19	Positif	0.99934	Positif	43	Positif	0.96592	Positif	67	Positif	0.96480	Positif
20	Positif	0.96480	Positif	44	Positif	0.96480	Positif	68	Positif	0.99934	Positif
21	Positif	0.96480	Positif	45	Positif	0.93273	Positif	69	Positif	0.84546	Positif
22	Positif	0.96480	Positif	46	Positif	0.96480	Positif	70	Positif	0.99934	Positif
23	Positif	0.09305	Negatif	47	Positif	0.98581	Positif	71	Positif	0.93273	Positif
24	Positif	0.96480	Positif	48	Positif	0.96480	Positif	72	Positif	0.96480	Positif















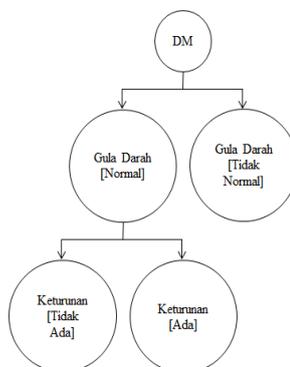








4.8.



Gambar 4.8 Pohon Hasil Iterasi Kedua

Selanjutnya dilakukan tahap rekursif pada cabang kiri. Data yang digunakan pada iterasi ketiga ini berdasarkan data cabang kiri dari iterasi kedua yang diperoleh *Gini Indeks* = 0,20245.

Kemudian, akan dilakukan perhitungan *Gini Indeks* juga pada masing-masing variabel X seperti pada langkah sebelumnya. Dari perhitungan tersebut, kriteria pemilihan pemilah iterasi ketiga disajikan pada Tabel 4.10.

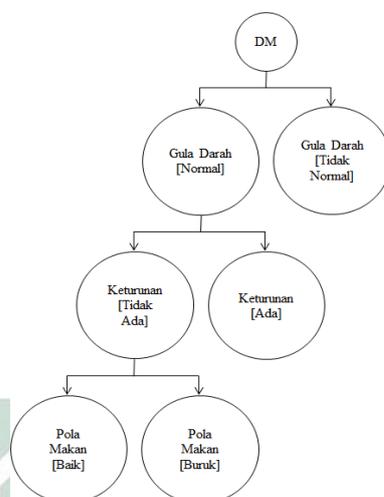
Tabel 4.10 Kriteria Pemilihan Pemilah Iterasi Ketiga

X	L				R				Gini A	Δ Gini
	Record	Negatif	Positif	Gini	Record	Negatif	Positif	Gini		
Usia	27	24	3	0,19753	8	7	1	0,21875	0,20238	6,8E-05
Keturunan	35	31	4	0,20245	0	0	0	0	0,20245	0
Gula Darah	35	31	4	0,20245	0	0	0	0	0,20245	0
Obesitas	9	9	0	0	26	22	4	0,26036	0,19341	0,00904
<b>Pola Makan</b>	<b>7</b>	<b>4</b>	<b>3</b>	<b>0,4898</b>	<b>28</b>	<b>27</b>	<b>1</b>	<b>0,06888</b>	<b>0,15306</b>	<b>0,04939</b>

Berdasarkan Tabel 4.10, dari kelima variabel tersebut nilai  $\Delta Gini X_5$  paling besar sehingga variabel Pola Makan dipilih menjadi pemilah ketiga. Pasien



## 4.9.



Gambar 4.9 Pohon Hasil Iterasi Ketiga

Selanjutnya dilakukan tahap rekursif pada cabang kiri. Data yang digunakan pada iterasi keempat ini berdasarkan data cabang kiri dari iterasi ketiga yang diperoleh *Gini Indeks* = 0,4898.

Kemudian, akan dilakukan perhitungan *Gini Indeks* juga pada masing-masing variabel X seperti langkah sebelumnya. Dari perhitungan tersebut, kriteria pemilihan pemilah iterasi keempat disajikan pada Tabel 4.11.

Tabel 4.11 Kriteria Pemilihan Pemilah Iterasi Keempat

X	L				R				Gini A	Δ Gini
	Record	Negatif	Positif	Gini	Record	Negatif	Positif	Gini		
Usia	4	2	2	0,5	3	2	1	0,44444	0,47619	0,01361
Keturunan	7	4	3	0,4898	0	0	0	0	0,4898	0
Gula Darah	7	4	3	0,4898	0	0	0	0	0,4898	0
Obesitas	2	2	0	0	5	2	3	0,48	0,34286	0,14694
Pola Makan	7	4	3	0,4898	0	0	0	0	0,4898	0





Tabel 4.12 Kriteria Pemilihan Pemilah Iterasi Kelima

X	L				R				Gini A	$\Delta$ Gini
	Record	Negatif	Positif	Gini	Record	Negatif	Positif	Gini		
Usia	2	2	0	0	0	0	0	0	0	0
Keturunan	2	2	0	0	0	0	0	0	0	0
Gula Darah	2	2	0	0	0	0	0	0	0	0
Obesitas	2	2	0	0	0	0	0	0	0	0
Pola Makan	2	2	0	0	0	0	0	0	0	0

Berdasarkan Tabel 4.12, dari kelima variabel tersebut nilai  $\Delta Gini = 0$  sehingga proses rekursif pada cabang kiri berhenti dan ditetapkan sebagai *class* Negatif, dikarenakan pada iterasi kelima data hanya terdiri dari *class* Negatif. Dengan demikian, pohon yang terbentuk dari iterasi kelima seperti Gambar 4.11.











Tabel 4.15 Kriteria Pemilihan Pemilah Iterasi Kedelapan

X	L				R				Gini A	$\Delta$ Gini
	Record	Negatif	Positif	Gini	Record	Negatif	Positif	Gini		
Usia	0	0	0	0	3	2	1	0,44444	0,44444	0
Keturunan	3	2	1	0,44444	0	0	0	0	0,44444	0
Gula Darah	3	2	1	0,44444	0	0	0	0	0,44444	0
Obesitas	0	0	0	0	3	2	1	0,44444	0,44444	0
Pola Makan	3	2	1	0,44444	0	0	0	0	0,44444	0

Berdasarkan Tabel 4.15, dari kelima variabel tersebut nilai  $\Delta Gini = 0$  sehingga proses rekursif pada cabang kanan berhenti dan ditetapkan sebagai *class* Negatif, dikarenakan pada iterasi kedelapan data *record* < 5 serta peluang *class* Negatif lebih besar dari pada *class* Positif yaitu 0,67 : 0,33. Dengan demikian, pohon yang terbentuk dari iterasi kedelapan seperti Gambar 4.14.



































Apabila pasien mempunyai kadar gula darah yang tidak normal, mempunyai tingkat kegemukan (obesitas) tidak normal dan berusia tua, maka pasien tersebut dikatakan positif Diabetes Mellitus.

Dalam penelitian ini, perbandingan data *training* dan data *testing* yang digunakan yaitu 70% : 30%. Berdasarkan data *training* tersebut terbentuk pohon klasifikasi optimal dengan ketepatan klasifikasi (akurasi) untuk data *training* sebesar 90,48% dan data *testing* sebesar 90,27%.

Dengan demikian, berdasarkan model yang terbentuk dari masing-masing metode dapat disimpulkan bahwa untuk data *training* metode Regresi Logistik Biner mempunyai tingkat akurasi lebih tinggi 2,38% daripada CART. Sementara itu, untuk data *testing* metode Regresi Logistik Biner mempunyai tingkat akurasi lebih tinggi 1,4% daripada CART. Oleh karena itu, pada penelitian ini metode Regresi Logistik Biner lebih baik dibandingkan metode CART.



3. Ketepatan klasifikasi dari model Regresi Logistik Biner sebesar 92,86% untuk data *training* dan data *testing* sebesar 91,67%. Sementara itu, akurasi dari pohon klasifikasi optimal CART sebesar 90,48% untuk data *training* dan data *testing* sebesar 90,27%. Dari kedua model diperoleh, untuk kasus klasifikasi penyakit Diabetes Mellitus dengan 5 variabel independen yaitu usia, faktor keturunan, kadar gula darah, obesitas dan pola makan metode Regresi Logistik Biner lebih baik apabila dibandingkan dengan metode CART. Hal ini dikarenakan metode Regresi Logistik Biner mempunyai tingkat akurasi lebih tinggi dibandingkan metode CART baik pada *training* maupun data *testing*.

## 5.2. Saran

Berdasarkan analisis hasil yang telah dilakukan, adapun saran yang dapat diberikan adalah sebagai berikut:

1. Dalam kasus klasifikasi penyakit Diabetes Mellitus menggunakan Regresi Logistik Biner maupun CART belum 100% akurat, sehingga untuk penelitian selanjutnya dapat menggunakan metode lain, misalnya *Random Forest*, SVM atau metode lainnya agar diperoleh nilai akurasi yang lebih tinggi.
2. Untuk penelitian berikutnya, dapat dilakukan penambahan faktor-faktor yang mempengaruhi penyakit Diabetes Mellitus, misalnya jenis kelamin, kadar insulin, tekanan darah dan faktor yang lainnya.
3. Jumlah data yang digunakan dapat diperbanyak, semakin banyak data yang digunakan semakin kecil nilai *error*nya.





- Diabetes UK. 2010. *Diabetes in the UK 2010. Key Statistics on Diabetes*: United Kingdom.
- Dorland, W. 2005. *Kamus Kedokteran Dorland*. EGC: Jakarta.
- Ghozali, Imam. 2005. *Aplikasi Analisis Multivariate Dengan Program SPSS*. Penerbit Universitas Diponegoro: Semarang.
- Goronescu, F. 2011. *Data Mining: Concept, Models and Techniques*. Springer: Verlag Berlin Heidelberg.
- Hosmer, D. W., dan Lemeshow. S. 2000. *Applied Logistic Regression*. John Wiley & Sons, Inc: New York.
- International of Diabetic Ferderation*. 2017. *IDF Diabetes Atlas*. Eighth Edition. [www.idf.org/diabetesatlas](http://www.idf.org/diabetesatlas) (Diakses tanggal 7 Oktober 2019)
- Johnson, R.A. dan Wichern, D. W. 2007. *Applied Multivariate Statistical Analysis*. Sixth Edition. Pearson Education Inc: USA.
- Kaban, Sempakata. 2007. *Diabetes Tipe 2 di Kota Sibolga Tahun 2005*. Majalah Kedokteran Nusantara. Volume 40 No. 2. Medan.
- Kementerian Kesehatan. 2010. *Petunjuk Teknis Pengukuran Faktor Risiko Diabetes Mellitus*.
- Kusrini, Luthfi, E. T. 2009. *Algoritma Data Mining*. Penerbit Andi: Yogyakarta.
- Lestawati, R., Rais dan Utami, I. T. 2018. *Perbandingan antara Metode CART (Classification and Regression Tress) dan Regresi Logistik (Logistic Regression) dalam Mengklasifikasikan Pasien Penderita DBD (Demam Berdarah Dengue)*. Jurnal Ilmiah Matematika dan Terapan, Volume 15 No. 1, pp 98-107,

Program Studi Matematika, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Tadulako. Palu.

Lewis dan Roger, J. 2000. *An Introduction to Classification and Regression Trees (CART) Analysis*. Annual Meeting of Society for Academic Emergency Medicine of Sanfransisco, UCLA Medical Center: California.

Margasari, A. 2014. *Penerapan Metode CART (Classification And Regression Trees) dan Analisis Regresi Logistik Biner pada Klasifikasi Profil Mahasiswa FMIPA Universitas Brawijaya*. Jurnal Mahasiswa Statistik, Vol. 2, No. 4, FMIPA Universitas Brawijaya. Malang.

Misnadiarly. 2006. *Diabetes Mellitus Gangren, Ulcer, Infeksi, Mengenali Gejala, Menanggulangi dan Mencegah Komplikasi*. Pustaka Obor Populer: Jakarta.

Mustamir. 2008. *5 Metode Penyembuhan dari Langit*. Lingkaran: Yogyakarta.

Rizki, F., Widodo, D. A., Wulandari, S. P. 2015. *Faktor Risiko Penyakit Anemia Gizi Besi pada Ibu Hamil di Jawa Timur Menggunakan Analisis Regresi Logistik*. Jurnal Sains dan Seni ITS. Vol. 4, No. 2;2337 - 3520.

Rumaendra, Wella. 2016. *Perbandingan Klasifikasi Penyakit Hipertensi Menggunakan Regresi Logistik Biner dan Algoritma C4.5*. Universitas Diponegoro: Semarang.

Sudoyo, A., Setiyohadi, B., Simadibrata, M., dan Setiati, S. 2009. *Buku Ajar Ilmu Penyakit dalam Jilid 3*. Edisi ., Interna Publishing: Jakarta.

WHO. 2006. *Definition, Diagnosis and Classification of Diabetes Mellitus and its Complication*. WHO.

