

**PENERAPAN METODE REGRESI BERSTRUKTUR POHON PADA
PENDUGAAN LAMA PENYUSUNAN SKRIPSI MAHASISWA**

ARTIKEL ILMIAH

**Artikel Ilmiah Ini Diambil Dari Sebagian Skripsi Untuk Memenuhi
Persyaratan Penyelesaian Program Sarjana Sains Jurusan Matematika
Fakultas Matematika dan Ilmu Pengetahuan Alam Universitas Jember**

Oleh :

THERESIA TRIAS CANDRA DEWI

NIM. 001810101012



**JURUSAN MATEMATIKA
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM
UNIVERSITAS JEMBER
FEBRUARI 2006**

PENGESAHAN

Artikel Ilmiah ini diterima oleh Fakultas Matematika dan Ilmu Pengetahuan Alam
Universitas Jember pada:

Hari :

Tanggal :

Tempat : Fakultas Matematika dan Ilmu Pengetahuan Alam

Dosen Pembimbing

Ketua
(Dosen Pembimbing Utama)
Anggota)

Sekretaris
(Dosen Pembimbing

Yuliani Setia Dewi, S.Si, M.Si
M.Si

NIP. 132 258 183

Agustina Pradjaningsih, S.Si,

NIP. 132 257 933

**Penerapan Metode Regresi Berstruktur Pohon Pada Pendugaan Lama
Penyusunan Skripsi Mahasiswa**

*(Application of Tree Regression Methods for forecasting the period of writing
student thesis)*

Theresia Trias Candra Dewi¹, Yuliani Setia Dewi², Agustina Pradjaningsih²

¹Mahasiswa Jurusan Matematika FMIPA Universitas Jember

²Staf Pengajar Matematika Jurusan Matematika FMIPA Universitas Jember

ABSTRACT

The tree regression is one of the regression methods that can be used to find out the influence of independent X variable to the dependent Y variable. The tree regression method analyzes data by doing step by step isolation. The pruning tree process is used in order to get the optimal tree size. Pruning tree is done based on cost complexity. This thesis is written in order to apply the tree regression method to the data of graduated student from Faculty of MIPA at 2001-2005 that used to find out the variable which has an influence to the period of writing student thesis. The research result shows that the period for student to write thesis is influenced by Cumulative Prestige Index and major.

Keywords : *tree regression, pruning tree, cost complexity*

ABSTRAK

Penerapan Regresi Berstruktur Pohon pada Pendugaan Lama Penyusunan Skripsi Mahasiswa, Theresia Trias Candra Dewi, 001810101012, Skripsi, Februari 2006, Jurusan Matematika, Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Jember

Regresi berstruktur pohon (*tree regression*) merupakan salah satu metode regresi yang dapat digunakan untuk mengetahui pengaruh variabel bebas X terhadap variabel tak bebas Y. Metode regresi berstruktur pohon menganalisa data dengan melakukan penyekatan secara bertahap. Proses pemangkasan pohon (*pruning tree*) diperlukan untuk mendapatkan ukuran pohon yang optimal. Pemangkasan pohon dilakukan berdasarkan biaya kompleksitas (*cost complexity*). Penulisan skripsi ini bertujuan untuk mengaplikasikan metode regresi pohon pada data lulusan mahasiswa Fakultas MIPA tahun 2001-2005 yaitu digunakan untuk mengetahui variabel yang berpengaruh pada lama penyusunan skripsi mahasiswa. Hasil penelitian menunjukkan lama penyusunan skripsi mahasiswa dipengaruhi oleh IPK dan jurusan.

Kata kunci : *tree regression, pruning tree, cost complexity.*

PENDAHULUAN

Salah satu tujuan analisis data dalam statistika adalah untuk mengetahui apakah ada hubungan antara dua variabel atau lebih dan memperkirakan besarnya efek kuantitatif dari perubahan variabel-variabel tersebut. Dalam keperluan analisis data, digunakan analisis regresi untuk mengetahui bentuk hubungan antara dua atau lebih variabel dimana terdapat sebuah variabel yang akan diramalkan atau disebut juga variabel tak bebas (*dependent variable*) yang dituliskan dalam Y dan satu atau lebih variabel yang digunakan untuk meramalkan atau disebut variabel bebas (*independent variable*) yang dituliskan dalam X .

Metode regresi pohon merupakan salah satu cara yang menarik dalam melakukan eksplorasi data dan mengambil kesimpulan dalam analisis. Perbedaan regresi berstruktur pohon ini dengan regresi yang biasa digunakan adalah pada regresi berstruktur pohon pendugaan respon dilakukan pada kelompok-kelompok pengamatan yang dibentuk berdasarkan variabel-variabel bebasnya, bukan untuk keseluruhan data.

Metode regresi berstruktur pohon ini menganalisa suatu gugus data dengan cara menyekatnya menjadi beberapa anak gugus (simpul) secara bertahap. Tahap pertama, seluruh data disekat menjadi dua anak gugus kemudian diperiksa kembali secara terpisah dan dibagi lagi berdasarkan penyekat lainnya, demikian seterusnya sampai tercapai kriteria berhenti. Anak gugus yang tidak bisa disekat dinamakan simpul terminal, sedangkan anak gugus yang masih dapat disekat dinamakan simpul dalam.

Permasalahan yang dibahas dalam penelitian ini adalah bagaimana mengaplikasikan metode Regresi Berstruktur pohon pada data kelulusan mahasiswa dan mengidentifikasi variabel yang berpengaruh terhadap lama penyusunan skripsi mahasiswa.

TINJAUAN PUSTAKA

Regresi Berstruktur Pohon

Regresi berstruktur pohon merupakan salah satu metode yang menggunakan kaidah pohon keputusan yang dibentuk melalui suatu algoritma penyekatan yang dilakukan secara bertahap. Metode regresi berstruktur pohon ini menganalisa suatu gugus data dengan cara menyekatnya menjadi beberapa anak gugus (simpul)

secara bertahap. Tahap pertama, seluruh data disekat menjadi dua anak gugus kemudian diperiksa kembali secara terpisah dan dibagi lagi berdasarkan penyekat lainnya, demikian seterusnya sampai tercapai kriteria berhenti. Anak gugus yang tidak bisa disekat dinamakan simpul terminal, sedangkan anak gugus yang masih dapat disekat dinamakan simpul dalam.

Algoritma Pohon Regresi

Dalam analisis regresi berstruktur pohon diperlukan empat langkah dasar yaitu (Roger L, 2000):

1. proses pembangunan pohon;
2. penghentian proses pembangunan pohon, dimana dalam proses ini akan ditemukan pohon dengan ukuran yang cukup besar (pohon maksimal);
3. proses pemangkasan (*pruning*) untuk mendapatkan pohon yang cukup sederhana;
4. pemilihan dan pembentukan pohon yang optimal.

Aturan Penyekatan

Pohon regresi dibentuk melalui penyekatan data pada tiap simpul ke dalam dua simpul anak. Aturannya sebagai berikut :

1. setiap penyekatan tergantung pada nilai yang hanya berasal dari satu variabel bebas;
2. untuk mengubah numerik X_j , penyekatan berasal dari pertanyaan “apakah $X_j \leq c$ ” untuk $c \in \mathfrak{R}$. Jadi jika ruang contohnya berukuran n dan terdapat sebanyak-banyaknya n nilai amatan yang berbeda pada peubah X_j , maka akan terdapat sebanyak-banyaknya $n-1$ split yang berbeda yang dibentuk oleh gugus pertanyaan (“apakah $X_j \leq c_i$ ”), dengan $i = 1, 2, 3, \dots, n-1$ dan c_i adalah nilai tengah antara dua nilai amatan peubah X_j yang berbeda dan berurutan;
3. untuk variabel bebas yang berkategori, pemilihan yang terjadi berasal dari semua kemungkinan pemilihan berdasarkan terbentuknya dua anak gugus yang saling lepas (*disjoint*). Jika X_j merupakan variabel kategori nominal bertaraf L , maka akan terdapat $2^{L-1} - 1$ sekatan yang mungkin, sedangkan jika X_j merupakan peubah kategori ordinal maka akan ada $L-1$ sekatan yang mungkin.

Tahap Penyekatan

Menurut Breiman *et al* (1993), untuk menyekat suatu simpul dilakukan proses sebagai berikut :

1. menentukan semua penyekat yang mungkin untuk setiap variabel bebas;
2. memilih sekat yang terbaik dari kumpulan sekat dua anak simpul, yaitu simpul kiri dan simpul kanan.

Penyekatan terbaik adalah penyekatan yang memaksimumkan ukuran pemisahan antara dua simpul anak tersebut. Jumlah Kuadrat Sisaan (JKS) digunakan sebagai kriteria kehomogenan di dalam masing-masing simpul. Misalkan simpul g berisi anak contoh $\{(x_n, y_n)\}$ dan $n(g)$ adalah banyaknya amatan pada simpul g , nilai respon dalam suatu simpul g tersebut dapat dihitung sebagai berikut :

$$\bar{y}(g) = \frac{1}{n(g)} \sum_{x_n \in g} y_n$$

maka jumlah kuadrat sisaan simpul g adalah :

$$JKS(g) = \sum_{x_n \in g} [y_n - \bar{y}(g)]^2$$

Misalkan s menyekat simpul g menjadi simpul kiri g_L dan simpul kanan g_R .

Kriteria jumlah kudrat terkecil $\Phi(s, g)$ adalah :

$$\Phi(s, g) = R(g) - R(g_L) - R(g_R)$$

dengan:

$R(g)$: jumlah kuadrat sisaan pada simpul g atau $JKS(g)$

$R(g_L)$: jumlah kuadrat sisaan pada simpul kiri g_L atau $JKS(g_L)$

$R(g_R)$: jumlah kuadrat sisaan pada simpul kanan g_R atau $JKS(g_R)$

Sekat terbaik s^* adalah sekat yang memenuhi kriteria $\Phi(s^*, g) = \max_{s^* \in \Omega} \Phi(s, g)$;

dengan Ω adalah himpunan semua sekat s yang mungkin pada simpul g ;

3. algoritma pembentukan struktur pohon dilakukan pada setiap variabel sampai dipenuhi aturan penghentian tertentu. Kriteria yang sering dijadikan aturan penghentian adalah N_{\min} banyaknya obyek pengamatan pada setiap simpul akhir;

4. menyusun tingkatan dari semua sekatan terbaik dalam setiap variabel berdasarkan penurunan jumlah kudrat terkecil. Hal ini berarti bahwa sekat yang dipilih untuk dijadikan penyekat utama adalah sekat yang mampu memberikan penurunan jumlah kuadrat sisaan terbesar.

Pemangkasan Pohon

Prosedur pemangkasan dilakukan berdasarkan suatu ukuran biaya kompleksitas (Breiman et al, 1993). Ukuran biaya kompleksitas dari subpohon G_{maks} (pohon berukuran besar atau pohon maksimal), yaitu ukuran biaya dari G , yang didefinisikan sebagai:

$$R_{\alpha}(G) = R(G) + \alpha|\tilde{G}|;$$

dengan :

$R_{\alpha}(G)$: biaya kompleksitas dari G

$|\tilde{G}|$: banyaknya anggota dari gugus simpul akhir \tilde{G}

$R(G)$: didefinisikan sebagai $R(G) = \sum_{g' \in \tilde{G}} R(g')$ dengan $R(g')$

adalah jumlah kuadrat sisaan pada suatu simpul akhir g'

α : parameter kompleksitas dengan $\alpha \geq 0$

Pohon yang dipangkas adalah pohon yang memenuhi kriteria biaya kompleksitas minimum.

Estimasi Penduga $R^{ts}(G)$ dan $R^{CV}(G)$

Dalam memilih pohon terbaik dari deretan pohon yang terbentuk pada proses pemangkasan digunakan suatu penduga yang dinamakan penduga jujur bagi $R(G)$ (Breiman *et al*, 1993). Ada dua penduga jujur bagi $R(G)$ yaitu penduga contoh uji $R^{ts}(G)$ dan penduga validasi silang $R^{CV}(G)$.

$R^{ts}(G_k)$ didefinisikan sebagai :

$$R^{ts}(G_k) = \frac{1}{n_2} \sum_{(x_i, y_i) \in L_2} [y_i - \hat{y}_k(x_i)]^2$$

dengan

$R^{ts}(G_k)$: penduga contoh uji bagi G_k

n_2 : ukuran dari *test sample* L_2

y_i : nilai respon pada amatan ke- i ; $i = 1, 2, \dots, n_2$

$\hat{y}_k(x_i)$: pendugaan respon dari amatan ke- i pada pohon ke- k

Pohon terbaik adalah G_{k_0} yang memenuhi : $R^{ts}(G_{k_0}) = \min R^{ts}(G_k)$.

Penduga validasi silang $R^{cv}(G_k)$ adalah sebagai berikut:

$$R^{cv}(G_k) = \frac{1}{N} \sum_{v=1}^v \sum_{(x_i, y_i) \in L_v} (y_i - \hat{y}_k^v(x_i))^2 ;$$

dengan N adalah jumlah amatan keseluruhan. Pohon terbaik adalah pohon G_{k_0} yang memenuhi kriteria:

$$R^{cv}(G_{k_0}) = \min R^{cv}(G_k)$$

METODOLOGI PENELITIAN

Pengumpulan Data

Data yang digunakan adalah data lulusan mahasiswa FMIPA Universitas Jember periode I bulan November 2001 sampai periode III bulan Maret 2005. Pengambilan data dilakukan pada bulan April 2005.

Identifikasi Variabel Penelitian

Pada penelitian ini variabel tak bebas (Y) adalah lama penyusunan skripsi mahasiswa FMIPA Universitas Jember, sedangkan variabel bebasnya adalah lulusan mahasiswa yang dibedakan atas:

1. Jurusan :
 - a. Matematika
 - b. Fisika
 - c. Biologi
 - d. Kimia
2. Jenis Kelamin
 - a. Laki-laki
 - b. Perempuan

3. Asal Daerah
 - a. Jember
 - b. Luar Jember
4. Jalur Masuk
 - a. PMDK
 - b. UMPTN
5. Indeks Prestasi Kumulatif (IPK)

HASIL DAN PEMBAHASAN

Gambaran Data

Data lulusan mahasiswa FMIPA tahun 2001-2005 adalah sebagai berikut :

	Jenis Kelamin		Asal Daerah		Jalur Masuk	
	L	P	Jember	Luar Jember	PMDK	UMPTN
Matematika	24	44	28	40	14	54
Fisika	20	45	29	36	17	48
Kimia	12	76	50	38	13	75
Biologi	25	51	36	40	9	67
Total	81	216	143	154	53	244

Sedangkan rata-rata IPK dan lama penyusunan skripsi mahasiswa

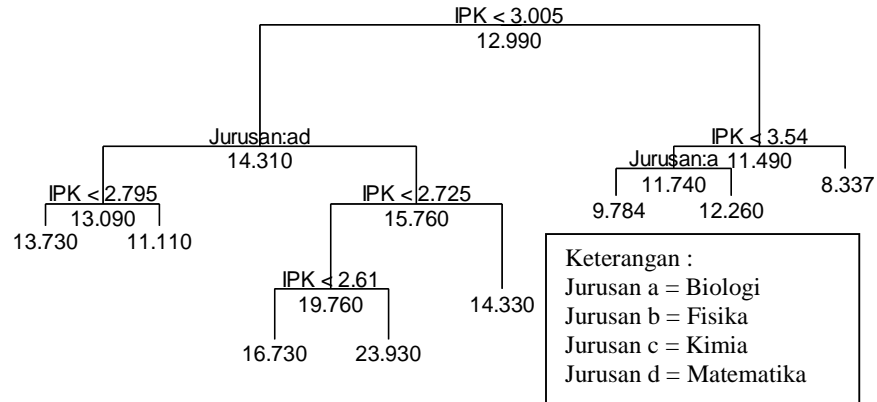
	Rata-rata IPK	Rata-rata Lama Penyusunan Skripsi
Matematika	3.04	11.32
Fisika	2.92	11.43
Kimia	2.79	12.77
Biologi	3.04	13.81
FMIPA	2.94	12.99

Pembangunan Pohon Regresi

Formula yang digunakan dalam pembangunan pohon regresi dengan menggunakan paket **R** adalah :

```
mipa.tree<-tree(Lama~JK+Jurusan+Asal+JlrMasuk+IPK,MIPA)
```

Regresi Pohon Fakultas MIPA



Gambar 1. Pohon Regresi Awal

Berdasarkan proses pembangunan pohon dan hasil plot pohon regresi terlihat bahwa terdapat 8 simpul akhir pada pembangunan pohon awal yaitu :

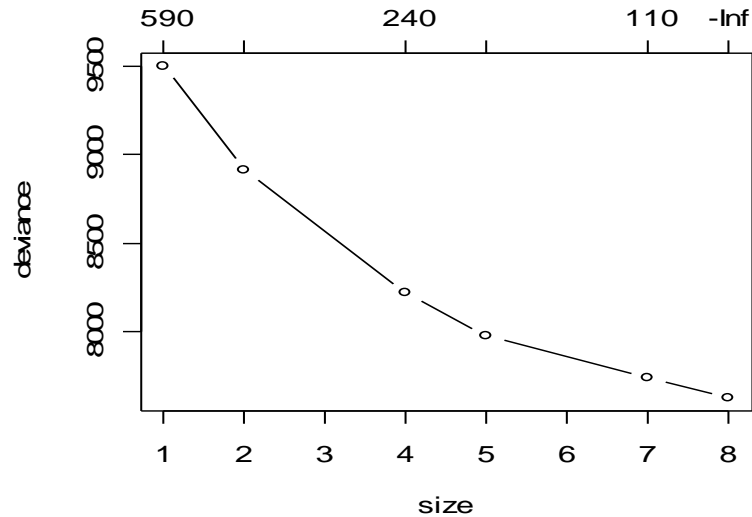
1. kelompok amatan yang memiliki $IPK < 2.795$ dari jurusan Matematika dan Biologi dengan rata-rata lama penyusunan skripsi 13.73 bulan ;
2. kelompok amatan yang memiliki $2.795 < IPK < 3.005$ dari jurusan Matematika dan Biologi dengan rata-rata lama penyusunan skripsi 11.11 bulan ;
3. kelompok amatan yang memiliki $IPK < 2.61$ dari jurusan Fisika dan Kimia dengan rata-rata lama penyusunan skripsi 16.73 bulan ;
4. kelompok amatan yang memiliki $2.61 < IPK < 2.725$ dari jurusan Fisika dan Kimia dengan rata-rata lama penyusunan skripsi 23.93 bulan ;
5. kelompok amatan yang memiliki $2.725 < IPK < 3.005$ dari jurusan Fisika dan Kimia dengan rata-rata lama penyusunan skripsi 14.33 bulan ;
6. kelompok amatan yang memiliki $3.005 < IPK < 3.54$ dari jurusan Biologi dengan rata-rata lama penyusunan skripsi 9.78 bulan ;
7. kelompok amatan yang memiliki $3.005 < IPK < 3.54$ dari jurusan Matematika, Fisika, dan Kimia dengan rata-rata lama penyusunan skripsi 12.26 bulan ;
8. kelompok amatan yang memiliki $IPK > 3.54$ dengan rata-rata lama penyusunan skripsi 8.34 bulan.

Pemangkasan Pohon

Pemangkasan pohon dilakukan dengan proses *pruning*. Perintah yang dilakukan pada paket **R** :

```
> prune.mipa<-prune.tree(mipa.tree)
```

Gambar- 4.2 Hubungan ukuran pohon dan deviansi dalam pemangkasan

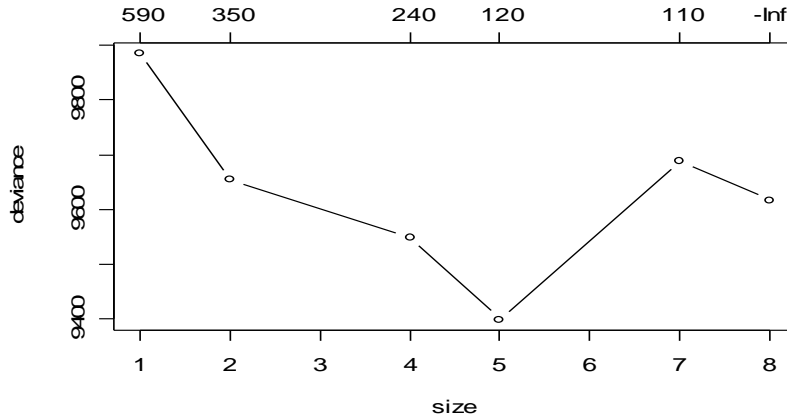


Plot diatas menunjukkan bahwa semakin banyak jumlah pohon yang terbentuk semakin kecil deviansinya. Hal ini menunjukkan bahwa biaya kompleksitas (*cost-validation*) yang dibutuhkan untuk pemangkasan akan semakin kecil.

Dalam menentukan ukuran pohon digunakan validasi silang (*cross-validation*). Formula yang digunakan dalam paket R untuk melakukan validasi silang (*cross validation*) adalah `cv.tree`.

```
> cv.mipa<-cv.tree(mipa.tree)
```

Hasil plot hubungan deretan pohon dan deviansi berdasarkan validasi silang (*cross validation*) adalah sebagai berikut :

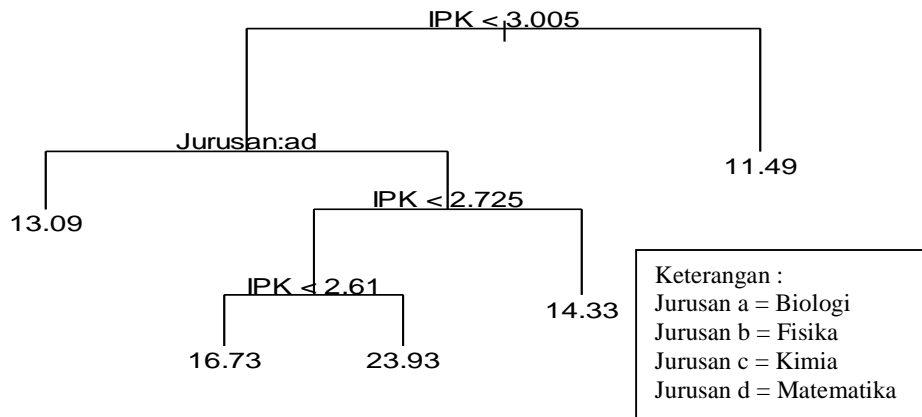


Gambar 3 Hubungan ukuran pohon dan deviansi berdasarkan validasi silang

Dari hasil plot `cv.tree` tampak bahwa ukuran pohon yang paling optimal dalam meminimumkan deviansi adalah simpul dengan ukuran (*size*) 5 yang menunjukkan bahwa pohon optimal memiliki 5 simpul akhir.

Pembentukan pohon optimal yang berukuran 5 adalah sebagai berikut :

```
> prune.mipa<-prune.tree(mipa.tree,best=5)
```



Gambar 4. Plot pohon optimal

Berdasarkan kelompok yang dihasilkan dengan menggunakan analisa regresi pohon diatas dapat disusun urutan kelompok mahasiswa berdasarkan waktu penyusunan skripsi yang lebih cepat sebagai berikut :

1. kelompok mahasiswa dengan waktu penyusunan skripsi paling cepat berada pada simpul 3 yaitu mahasiswa yang memiliki $IPK > 3,005$ dengan rata-rata lama penyusunan skripsi 11,49 bulan;
2. kelompok mahasiswa dengan waktu penyusunan skripsi cepat berada pada simpul 4 yaitu mahasiswa jurusan Matematika dan mahasiswa jurusan Biologi yang memiliki $IPK < 3,005$ dengan rata-rata lama penyusunan skripsi 13,09 bulan;
3. kelompok mahasiswa dengan waktu penyusunan skripsi sedang berada pada simpul 11 yaitu mahasiswa jurusan Fisika dan mahasiswa jurusan Kimia yang memiliki $2,725 < IPK < 3,005$ dengan rata-rata lama penyusunan skripsi 14,33 bulan;
4. kelompok mahasiswa dengan waktu penyusunan skripsi lama berada pada simpul 20 yaitu mahasiswa jurusan Fisika dan mahasiswa jurusan Kimia yang memiliki $IPK < 2,61$ dengan rata-rata lama penyusunan skripsi 16,73 bulan;
5. kelompok mahasiswa dengan waktu penyusunan skripsi paling lama berada pada simpul 21 yaitu mahasiswa jurusan Fisika dan mahasiswa jurusan Kimia yang memiliki $2,61 < IPK < 2,725$ dengan rata-rata lama penyusunan skripsi 23,93 bulan.

KESIMPULAN DAN SARAN

Kesimpulan

Dari hasil analisis data yang dibahas pada Bab IV dapat diambil kesimpulan sebagai berikut :

1. hasil analisa data lulusan mahasiswa FMIPA tahun 2001-2005 menunjukkan bahwa variabel yang paling berpengaruh terhadap lama penyusunan skripsi mahasiswa adalah variabel IPK dan Jurusan ;
2. hasil analisa regresi pohon pada data lulusan mahasiswa FMIPA tahun 2001-2005 berdasarkan IPK, menunjukkan bahwa waktu penyusunan skripsi yang paling cepat terdapat pada kelompok mahasiswa dengan $IPK > 3,005$ dengan rata-rata lama penyusunan skripsi 11,49 bulan ;

3. hasil analisa regresi pohon pada data lulusan mahasiswa FMIPA tahun 2001-2005 berdasarkan jurusan dari kelompok mahasiswa yang memiliki $IPK < 3,005$ menunjukkan mahasiswa jurusan Matematika dan Biologi memiliki waktu penyelesaian skripsi yang lebih cepat dibandingkan dengan mahasiswa jurusan Fisika dan Kimia.

Saran

Diharapkan agar dapat menggunakan regresi pohon ini dalam analisis dan dapat dikembangkan dengan menggunakan software analisis yang lain seperti Paket S-Plus, SAS dan sebagainya. Dari hasil penelitian ini diharapkan ada penelitian selanjutnya yang dapat menggambarkan kelompok variabel-variabel lainnya yang dalam penelitian ini tidak digunakan dalam membangun pohon regresi.

DAFTAR PUSTAKA

- Breiman, L. J.H Friedman, R. A. Olshen & Charles J. Stone. 1993. *Classification and Regression Tree*. Chapman & Hall. New York
- Chambers J.M & Hastie T.J. 1993. *Statistica Model in S*. Chapman & Hall. New York
- Lewis, Roger J. 2000. *An Introduction to Classification and Regression Tree (CART) Analysis*. Department of Emergency Medicine Harbor-UCLA Medical Centre. California. <http://www.saem.org/download/lewis1.pdf>
- Venables W. N & Ripley B.D. 1994. *Modern Applied Statistics with S-Plus*. Springer. New York
- Yohannes Y & Hoddinot J. 1999. *Classification and Regression Trees : An Introduction*. International Food Policy Research Institute. USA

DAFTAR ISI

HALAMAN JUDUL.....	i
HALAMAN MOTTO.....	ii
HALAMAN PERSEMBAHAN.....	iii
HALAMAN DEKLARASI.....	iv
ABSTRAK.....	v
HALAMAN PENGESAHAN.....	vi
KATA PENGANTAR.....	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	x
DAFTAR GAMBAR.....	xi
DAFTAR LAMPIRAN.....	xii
I. PENDAHULUAN	
1.1 Latar Belakang.....	1
1.2 Permasalahan.....	2
1.3 Tujuan.....	3
1.4 Manfaat.....	3
II. TINJAUAN PUSTAKA	
2.1 Regresi Pohon.....	4
2.2 Algoritma Pohon Regresi.....	6
2.3 Aturan Penyekatan.....	7
2.4 Tahap Penyekatan.....	7
2.5 Penentuan Ukuran Pohon.....	8
2.5.1 Pemangkasan Pohon (<i>Prunning Tree</i>).....	9
2.5.2 Estimasi Penduga $R^{ts}(G)$ dan $R^{cv}(G)$	11
2.6 Paket R.....	13
III. METODOLOGI PENELITIAN	
3.1 Pengumpulan Data.....	15
3.2 Identifikasi Obyek Populasi, dan Sampel Penelitian.....	15
3.3 Identifikasi Variabel Penelitian.....	15

3.4 Metode Pengolahan Data.....	16
IV.HASIL DAN PEMBAHASAN	
4.1 Gambaran Data.....	17
4.2 Pembangunan Pohon Regresi.....	18
4.3 Pemangkasan Pohon.....	21
4.4 Dugaan Lama Studi Mahasiswa.....	24
V. KESIMPULAN DAN SARAN	
5.1 Kesimpulan.....	28
5.2 Saran.....	28
DAFTAR PUSTAKA.....	29
LAMPIRAN.....	30