

METODE KLASIFIKASI BERSTRUKTUR POHON BINER (STUDI KASUS PADA PRAKIRAAN SIFAT HUJAN BULANAN DI BOGOR)^{1)T}

Aan Kardiana²⁾, Aunuddin³⁾, Aji Hamim Wigena³⁾, Hari Wijayanto³⁾

²⁾Jurusan Teknik Informatika, FTI, Universitas YARSI, Jakarta

E-mail: kardiana@yarsi.ac.id

³⁾Program Pascasarjana, Institut Pertanian Bogor

ABSTRAKSI

Makalah ini menguraikan tentang aplikasi pohon biner berupa pembentukan Pohon Klasifikasi dengan menggunakan algoritma pemilahan secara rekursif dilanjutkan dengan kombinasi strategi pemangkasan pohon ber-cost-complexity terkecil dan penggunaan contoh validasi silang, untuk mendapatkan pohon yang berukuran tepat. Metode ini diterapkan dalam prakiraan sifat hujan bulanan menggunakan data klimatologi rata-rata bulanan di stasiun Klimatologi Kelas I Darmaga Bogor. Peubah rata-rata suhu pada jam 07.00 wib. merupakan peubah dominan. Hasil evaluasi prakiraan selama periode pengujian, tingkat ketepatan metode: Regresi Linier Osilasi Selatan; Probabilitas; Regresi Linier Curah Hujan dan Pohon Klasifikasi masing-masing bernilai: 42,86 %; 14,29 %; 14,29 %; dan 71,43%.

Kata kunci: Pohon Klasifikasi.

1. PENDAHULUAN

Salah satu aplikasi dari pohon biner adalah metode Pohon Klasifikasi. Penggunaan model berbasis pohon ini masih sedikit, meskipun berbagai penelitian telah menunjukkan bahwa pohon ini merupakan salah satu cara menarik dalam melakukan eksplorasi data dan pengambilan kesimpulan. Apalagi saat ini telah didukung oleh teknologi komputer yang berkembang semakin pesat sehingga proses komputasi yang melibatkan perhitungan yang cukup rumit dan mengandung banyak iterasi dapat dilakukan dengan semakin mudah dan cepat sehingga memberikan hasil yang lebih: cepat; tepat dan akurat.

Perancangan pengambilan keputusan berbasis pohon klasifikasi bermula dari studi dalam bidang ilmu sosial pada tahun 1973 oleh Morgan & Messenger dengan menggunakan program *THAID* (*Theta Automatic Interaction Detection*). Dalam bidang statistika, Breiman *et al.* telah merancang pengambilan keputusan berbasis pohon dengan mengembangkan konsep-konsep sehingga berpengaruh pada dua hal yaitu membawa masalah untuk diperhatikan oleh statistikawan dan mengusulkan algoritma untuk membentuk pohon [3]. Pengantar dalam S tentang model berbasis pohon telah mengembangkan metode dengan fungsi-fungsi untuk memeriksa dan mengevaluasi pohon serta mempermudah penggunaannya sehingga banyak memberikan dukungan untuk eksplorasi dan analisis data [4].

Penyusunan pohon klasifikasi dapat juga dilihat sebagai salah satu jenis pemilihan peubah dan untuk ukuran yang besar merupakan transformasi monotonik dari peubah penjelas x dan peubah respon y . Kita akan dapat mengetahui manakah peubah yang dominan di antara sederet peubah yang ada serta dapat pula mengidentifikasi

peubah-peubah yang hanya berpengaruh secara lokal dalam kelompok tertentu.

Metode ini telah diaplikasikan dalam bidang: kedokteran; pertahanan dan keamanan; fisika; dan lain-lain. Salah satu aplikasi di bidang teknologi informasi adalah dalam masalah pengenalan digit mesin hitung (kalkulator) elektrik.

Makalah ini bertujuan untuk menguraikan metode pohon klasifikasi dan aplikasinya dalam pembuatan prakiraan sifat hujan bulanan. Hasil yang diperoleh akan dibandingkan dengan hasil keluaran menurut metode yang sekarang digunakan oleh BMG yaitu: metode Probabilitas; metode Regresi Linier Osilasi Selatan; dan metode Regresi Linier Curah Hujan [6].

2. TEORI PENUNJANG

2.1 Pohon Klasifikasi

Misalkan (X, Y) merupakan peubah acak dengan $X=(X_1, X_2, \dots, X_q) \subset \mathcal{R}^q$ tersusun dari q vektor penjelas dan $Y \in \mathcal{Z} = \{1, 2, \dots, J\}$ vektor label kelas yang berkaitan, sehingga terdapat himpunan data [5] atau *Learning sample* [3] $L^{(0)} = \{(x_n, j_n), n = 1, 2, \dots, N^{(0)}\}$. Tujuan dari metode klasifikasi adalah ingin menaksir Y jika diberikan X .

Pohon klasifikasi terdiri atas $T = \{1, 2, 3, \dots\}$ dan fungsi $l(t)$ dan $r(t)$ dari T ke $T + \{0\}$ yang memenuhi ($l(t) = 0$ dan $r(t) = 0$) atau ($l(t) > t$ dan $r(t) > t$) serta $\forall t \in T - \{1\}, \exists! s \in T \ni t = l(s)$ atau $t = r(s)$ [5]. T dinamakan pohon klasifikasi; $\forall t \in T$ dinamakan simpul pohon; Elemen minimum dari T dinamakan *root*(T); Jika $s, t \in T$ dan $t = l(s)$ atau $t = r(s)$, maka s dinamakan *parent*(t) dan t dinamakan *sons*(s), *parent*(*root*(T)) = 0; $t \in T \ni$ *sons*(t) = 0, dinamakan simpul terminal sedangkan $t \neq 1 \in T - \{\text{sons}(t) = 0\}$ dinamakan simpul internal; Jika $s = \text{parent}(t)$ atau $s = \text{parent}(\text{parent}(t)) \dots$, t dinamakan *descendant*(s) dan s dinamakan *ancestor*(t);

$\forall t \in T$, cabang dari T dengan $root\ t \in T$ ditulis T_t , adalah himpunan bagian dari T yang memuat t dan seluruh *descendant* t ; Pemangkasan cabang T_t dari T dengan cara menghilangkan seluruh *descendant*(t) dari T yaitu potong seluruh T_t kecuali $root(T_t)$. Pohon terpangkas dinotasikan dengan $T-T_t$; Cabang dari T_t dari T dinotasikan dengan T_{t_1} ; T_t pohon bagian terpangkas dari T , dinotasikan dengan $T_{t_1} \leq T$ jika T_{t_1} pohon bagian dari T dan $root(T_{t_1}) = root(T)$. $T_{t_1} < T$ berarti $T_{t_1} \leq T$ dan $T_{t_1} \neq T$; \tilde{T} adalah himpunan terminal dari T [7].

Pohon klasifikasi dapat juga dipandang sebagai suatu aturan keputusan yang merupakan fungsi $d(\cdot)$ dari R^q ke ζ [5] atau $\forall x \in X, d(x) = y \in \zeta$ [3]. Jadi $d(\cdot)$ mempartisi ruang X menjadi J himpunan saling lepas $A(1), A(2), \dots, A(j)$ sehingga kaitan antara aturan keputusan dengan pohon klasifikasi dapat dinyatakan sebagai:

$$\forall x \in X \quad \forall t \in \tilde{T} \quad d(x) = j(t) \Leftrightarrow x \in A(t)$$

2.2 Pembentukan Pohon Klasifikasi

Pohon klasifikasi dibentuk melalui pemilahan secara iteratif terhadap $X = root(t) = A(1)$ kemudian dilanjutkan terhadap masing-masing $A(s), s \in sons(t), t \in T$, dan seterusnya.

Algoritma pembentukan pohon klasifikasi terdiri dari empat tahapan, yaitu: pemilihan pemilah; penentuan simpul terminal; penandaan label kelas; dan penentuan pohon dengan ukuran tepat.

2.2.1 Pemilihan Pemilah

Pada tahap ini dicari pemilah dari setiap simpul yang menghasilkan penurunan tingkat keheterogenan paling tinggi.

Keheterogenan suatu simpul diukur berdasarkan nilai *impurity*-nya. Nilai ini akan maksimum jika seluruh kelas secara seimbang bercampur di dalamnya dan akan minimum ketika hanya memuat satu kelas. Beberapa fungsi *impurity* yang dapat dipergunakan adalah indeks *Gini*: $f(p) = p(1-p)$ [3], indeks informasi: $f(p) = -\ln(p)$ dan indeks *twoing*: $I(t) = \min_{C_1, C_2} [f(p_{C_1}) + f(p_{C_2})]$ C_1, C_2

merupakan beberapa partisi dari J kelas menjadi dua himpunan bagian saling lepas. Jika $J=2$, *twoing* ekuivalen dengan indeks *impurity* biasa [9], atau indeks *entropi* $f(p) = N_j(t) \ln(p)$ [11].

Pemilahan dimulai dari ruang umum dengan cara memeriksa nilai-nilai dari setiap peubah penjelas. Peubah kontinu atau ordinal x_j memilah berbentuk $x_j \leq t$ lawan $x_j > t, t \in \mathcal{R}$ sedangkan untuk peubah nominal, L taraf dibagi menjadi dua himpunan saling lepas (akan terdapat sebanyak 2^{L-1} kemungkinan pemilah).

Untuk suatu t , misalkan terdapat calon pemilah s yang memilah t menjadi t_L (dengan proporsi p_L) dan menjadi t_R (dengan proporsi p_R),

maka kebaikan dari s didefinisikan sebagai penurunan *impurity*:

$$\Delta i(s, t) = i(t) - p_L i(t_L) - p_R i(t_R)$$

Pemilah s mengirimkan semua x_n di dalam t yang menjawab “ya” ke t_L dan semua x_n di dalam t yang menjawab “tidak” ke t_R .

Pohon dikembangkan dengan cara: pada simpul t_1 , carilah s^* yang memberikan nilai penurunan *impurity* tertinggi yaitu:

$$\Delta i(s^*, t_1) = \max_{s \in S} \Delta i(s, t_1)$$

maka t_1 dipilah menjadi t_2 dan t_3 menggunakan s^* . Dengan cara yang sama dilakukan juga pencarian pemilah terbaik pada t_2 dan t_3 secara terpisah, dan seterusnya.

2.2.2 Penentuan Simpul Terminal

Ketika suatu simpul t dicapai sehingga tidak terdapat penurunan *impurity* secara berarti ($\frac{\Delta i}{i} < 1\%$ devian *root*) [10] atau banyaknya objek cukup kecil ($n_i < 5$) [4] maka selanjutnya t tidak dipilah lagi tetapi dijadikan simpul terminal dan hentikan pembentukan pohon.

2.2.3 Penandaan Label Kelas

Misalkan bahwa T telah dikonstruksi menggunakan data dalam $L^{(0)}$. Pada masing-masing $t \in T$, nyatakanlah $p(t) = N(t)/N(0)$ [$N(t)$ = jumlah data anggota $L^{(0)}$ yang memenuhi $x_n \in A(t)$]; $p(t)$ merupakan taksiran dari peluang $P(x_n \in A(t))$ berdasar $L^{(0)}$. $P(Y=j|X \in A(t))$ berdasar $L^{(0)}$ ditaksir oleh $p(j|t) = N_j(t)/N(t), j = 1, 2, \dots, J, [N_j(t)$ adalah jumlah data anggota $L^{(0)}$ yang memenuhi $x_n \in A(t)$ dan $y_n = j]$.

Label kelas dari simpul terminal ditentukan berdasarkan aturan jumlah terbanyak, yaitu jika

$$p(j_0|t) = \max_j p(j|t) = \max_j \frac{N_j(t)}{N(t)}, \text{ maka label}$$

kelas untuk terminal t adalah j_0 .

2.3 Penentuan Pohon optimum

Pohon klasifikasi hasil proses pembentukan tidak dibatasi ukurannya. Pohon klasifikasi terbesar memberikan nilai *misclassification error* paling kecil sehingga kita akan selalu cenderung memilih pohon ini untuk prakiraan. Tetapi, pohon ini cukup kompleks dalam menggambarkan struktur data.

Salah satu cara penentuan pohon klasifikasi optimum --tanpa mengorbankan kebaikan kecocokan melalui pengurangan simpul pohon (sehingga dicapai penghematan gambaran)-- dikenal dengan istilah pemangkasan (*pruning*). Pemangkas berturut-turut memangkas pohon bagian yang kurang penting. Tingkat kepentingan sebuah pohon bagian diukur berdasar ukuran *cost-complexity* [10] adalah:

$$D_\alpha(T_k) = D(T_k) + \alpha |\tilde{T}_k|$$

dimana:

- $D(T_k)$ = devian dari pohon bagian T_k ,
- \tilde{T}_k = himpunan simpul terminal pada T_k ,
- $|\tilde{T}_k|$ = banyak simpul terminal pada \tilde{T}_k
- α = parameter *cost-complexity*.

Hasil proses pemangkasan berupa sederet pohon klasifikasi T_k dan dengan menggunakan contoh validasi silang (*cross-validation sample*) dapat ditentukan pohon optimum T_{k_0} dengan

$$R^{cv}(T_{k_0}) = \min_k R^{cv}(T_k).$$

2.4 Metode prakiraan BMG

Metode prakiraan sifat hujan bulanan yang sekarang digunakan oleh BMG adalah:

2.4.1 Metode Probabilitas

Analisis yang dihasilkan dari metode Probabilitas berupa nilai prosentase dari tiga kategori sifat hujan yaitu di bawah normal (BN), normal (N), dan di atas normal (AN). Nilai prosentase terbesar dari ketiga sifat tersebut merupakan hasil yang akan dipakai dalam suatu prakiraan [6].

2.4.2 Metode Regresi Linier Osilasi Selatan

Peubah respon pada metode Regresi Linier Osilasi Selatan adalah jumlah curah hujan tiga bulan terakhir, dan indeks osilasi selatan pada tiga bulan sebelumnya sebagai peubah prediktor [6].

2.4.3 Metode Regresi Linier Curah Hujan

Metode Regresi Linier Curah Hujan melibatkan data curah hujan bulan yang akan diprakirakan sebagai peubah respon dan jumlah curah hujan dua bulan sebelumnya sebagai prediktor [6].

3. DATA DAN METODE

Data yang digunakan dalam penelitian ini adalah data klimatologi rata-rata bulanan yang diamati di stasiun Klimatologi Kelas I Darmaga Bogor, dari bulan Februari tahun 1980 sampai bulan April 1999 [1].

Peubah responnya adalah sifat hujan bulanan (**s**) yang berjenis kategorik. Nilainya adalah **BN** jika jumlah curah hujan bulanan $< 0,85 \cdot \text{normal curah hujan}$; **N**; jika $0,85 \cdot \text{normal curah hujan} \leq \text{jumlah curah hujan bulanan} \leq 1,15 \cdot \text{normal curah hujan}$ **AN** jika jumlah curah hujan bulanan $\geq 1,15 \cdot \text{normal curah hujan}$. Sedangkan peubah penjelasnya terdiri dari peubah rata-rata : suhu ($^{\circ}\text{C}$) pada jam : 07.⁰⁰ wib (**t07**); 13.⁰⁰ wib (**t13**); 18.⁰⁰ wib (**t18**), suhu terbesar (**tmax**); suhu terkecil (**tmin**), tingkat penyinaran matahari dalam % (**rm**), tekanan udara dalam mb (**tu**), kelembaban nisbi dalam % pada jam : 07.⁰⁰ wib (**k07**); 13.⁰⁰ wib

(**k13**); 18.⁰⁰ wib (**k18**); sifat hujan bulan sebelumnya (**s₀**); dan indeks osilasi selatan (**osi**). Selain itu digunakan juga data-data prakiraan dan evaluasi sifat hujan bulanan yang telah diterbitkan oleh BMG [2] untuk membandingkan tingkat ketepatan metode.

4. HASIL DAN PEMBAHASAN

Hasil prakiraan berdasarkan metode : Probabilitas; Regresi Linier Osilasi Selatan; dan Regresi Linier Curah Hujan masing-masing memberikan tingkat ketepatan sebesar: 14,29%; 42,86%, dan 14,29%. Hasil selengkapnya dapat dilihat pada Tabel 1, Tabel 2, dan Tabel 3.

Tabel 1. Prakiraan Sifat Hujan Bulanan menurut metode Probabilitas.

Tahun	No	Bulan	S	\hat{S}_p	Hasil
1998	1	Oktober	AN	BN	Salah
	2	Nopember	BN	AN	Salah
	3	Desember	BN	AN	Salah
1999	4	Januari	BN	AN	Salah
	5	Februari	BN	N	Salah
	6	Maret	BN	BN	Benar
	7	April	N	BN	Salah

dimana:

S \equiv sifat hujan sebenarnya

\hat{S}_p \equiv sifat hujan menurut metode probabilitas

Tabel 2. Prakiraan Sifat Hujan Bulanan menurut metode Regresi Linier Osilasi Selatan

Tahun	No	Bulan	S	\hat{S}_{os}	Hasil
1998	1	Oktober	AN	AN	Benar
	2	Nopember	BN	AN	Salah
	3	Desember	BN	BN	Benar
1999	4	Januari	BN	BN	Benar
	5	Februari	BN	AN	Salah
	6	Maret	BN	AN	Salah
	7	April	N	AN	Salah

dimana:

\hat{S}_{os} \equiv sifat hujan menurut metode Regresi Linier Osilasi Selatan

Tabel 3. Prakiraan Sifat Hujan Bulanan menurut metode Regresi Linier Curah Hujan

Tahun	No	Bulan	S	\hat{S}_{ch}	Hasil
1998	1	Oktober	AN	N	Salah
	2	Nopember	BN	N	Salah
	3	Desember	BN	N	Salah
1999	4	Januari	BN	N	Salah
	5	Februari	BN	AN	Salah
	6	Maret	BN	N	Salah
	7	April	N	N	Benar

dimana:

\hat{S}_{ch} \equiv sifat hujan menurut metode Regresi Linier Curah Hujan

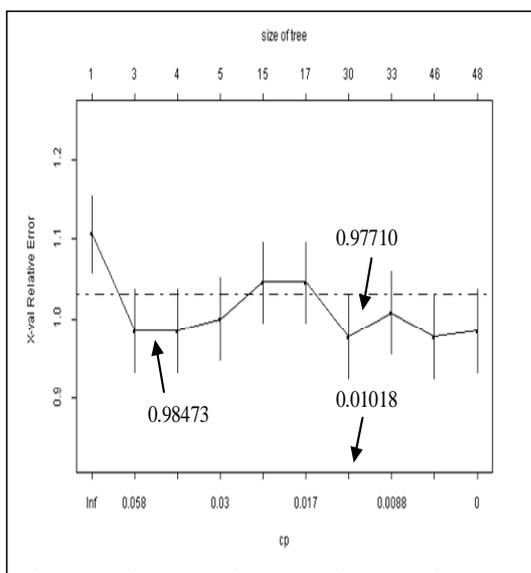
Penerapan metode Pohon klasifikasi pada data klimatologi rata-rata bulanan menghasilkan Pohon klasifikasi terbesar seperti terlihat pada Gambar 1.

Pohon klasifikasi yang telah diperoleh selanjutnya secara iteratif dipangkas menjadi deretan pohon bagian yang makin kecil dan tersarang.

Kemudian dengan mengkombinasikan aturan *cost complexity minimum* dan penggunaan contoh validasi silang *10-fold (10-fold cross-validation sample)*, diperoleh pohon berukuran tepat (optimum) yang dicapai pada parameter *complexity* bernilai $cp = 0.01018$ karena memberikan nilai R^{cv} (*x-val Relative Error*) terkecil yaitu $R^{cv} = 0.9771$ (Gambar 2).



Gambar 1. Pohon klasifikasi terbesar



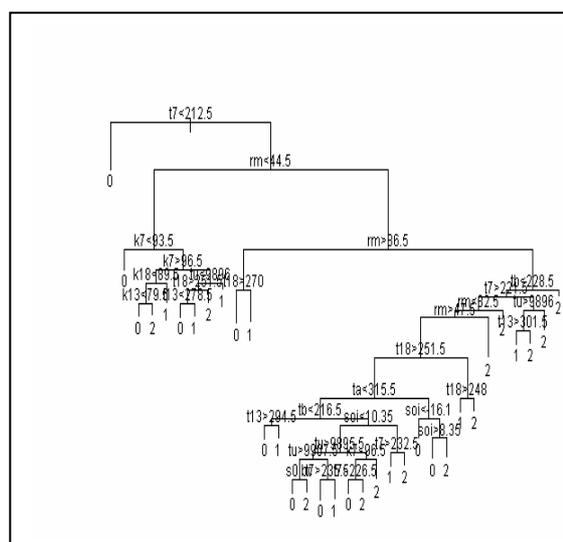
Gambar 2. Plot parameter *complexity*

Pohon optimum mengandung 30 buah simpul terminal dengan tingkat ketepatan pengelompokkan amatan sebesar 69.47% (Gambar 3).

Peubah penjelas yang pertama memilah amatan adalah peubah suhu pada jam 07.⁰⁰ wib. sehingga peubah ini merupakan peubah dominan. Hasil ini sesuai dengan Mc. Hattie & Schnelle [8] bahwa yang berperan penting dalam menentukan

klasifikasi sifat hujan adalah curah hujan dan suhu udara.

Amatan sebanyak 205 buah, pada simpul pertama (*root*) dipilah menjadi kelompok kiri atau kelompok kanan berdasarkan nilai peubah rata-rata suhu pada jam 07.⁰⁰ wib. Amatan yang memiliki nilai suhu pada jam 07.⁰⁰ wib < 21,25°C sebanyak 7 buah mengelompok di simpul 2 (simpul terminal) sedangkan amatan yang memiliki nilai suhu pada jam 07.⁰⁰ wib ≥ 21,25°C sebanyak 198 buah mengelompok di simpul 3. Kelompok amatan pada simpul 3 selanjutnya oleh peubah tingkat penyinaran matahari dipilah menjadi kelompok di simpul 6 yang terdiri atas 43 buah jika tingkat penyinaran matahari < 44,5% dan amatan dengan tingkat penyinaran matahari ≥ 44,5% sebanyak 155 buah mengelompok di simpul 7.



Gambar 3. Pohon klasifikasi optimum

Amatan pada simpul 6 dipilah lagi berdasar peubah kelembaban pada jam 07.⁰⁰ wib. Sebanyak 5 buah amatan yang mempunyai nilai peubah kelembaban pada jam 07.⁰⁰ wib < 93.5% mengelompok di simpul 12 (simpul terminal) sedangkan amatan lainnya dengan kelembaban pada jam 07.⁰⁰ wib ≥ 93.5% sebanyak 38 buah mengelompok di simpul 13. Proses ini dilanjutkan hingga dicapai semua simpul terminal.

Tabel 4. Prakiraan Sifat Hujan Bulanan menurut metode Pohon Klasifikasi

Tahun	No	Bulan	s	\hat{s}_o	Hasil
1998	1	Oktober	AN	AN	Benar
	2	Nopember	BN	N	Salah
	3	Desember	BN	BN	Benar
1999	4	Januari	BN	BN	Benar
	5	Februari	BN	N	Salah
	6	Maret	BN	BN	Benar
	7	April	N	N	Benar

dimana:

\hat{s}_o ≡ sifat hujan menurut pohon optimum

Selanjutnya pohon klasifikasi yang diperoleh digunakan untuk membuat prakiraan sifat hujan bulanan dengan cara menelusuri pohon klasifikasi optimum dari mulai simpul *root* hingga dicapai simpul terminal berdasarkan nilai peubah penjelas yang diperoleh.

Metode pohon klasifikasi optimum (yang mengandung 30 buah simpul terminal) memberikan tingkat ketepatan sebesar 71,43%. Hasil selengkapnya bisa dilihat pada Tabel 4.

KESIMPULAN

Kesimpulan yang diperoleh dalam penelitian ini adalah:

1. Peubah rata-rata suhu pada jam 07.⁰⁰ wib. merupakan peubah dominan.
2. Dari hasil evaluasi prakiraan selama 7 (tujuh) bulan terakhir, metode pohon klasifikasi optimum memberikan tingkat ketepatan sebesar 71,43 %.
3. Selama periode pengujian, dibandingkan dengan metode prakiraan yang sekarang digunakan oleh BMG, tingkat ketepatan metode pohon klasifikasi nilainya 28,57 % di atas tingkat ketepatan metode Regresi Linier Osilasi Selatan (yang bernilai 42,86 %), dan 57,14 % di atas tingkat ketepatan metode Probabilitas dan Regresi Linier Curah Hujan (yang keduanya bernilai 14,29 %). Jadi, metode pohon klasifikasi mempunyai tingkat ketepatan yang lebih baik dibandingkan dengan metode yang sekarang digunakan oleh BMG.

DAFTAR PUSTAKA

- [1] Badan Meteorologi dan Geofisika Balai Wilayah II. 1980-1999. *Data-data*

Klimatologi Bulanan. Stasiun Darmaga, Bogor.

- [2] Badan Meteorologi dan Geofisika Balai Wilayah II. 1998-1999. *Evaluasi dan Prakiraan Curah Hujan Propinsi Jawa Barat*. Stasiun Darmaga, Bogor.
- [3] Breiman, L., J. H. Friedman, R. A. Olshen & C.J. Stone. 1984. *Classification and Regression Tree*. Chapman and Hall, New York.
- [4] Clark, L.A. & D. Pregibon. 1992. *Tree-based models*. In Chambers, J.M. & T. J. Hastie (eds). 1992. *Statistical Models in S*. Chapman and Hall, New York.
- [5] Gelfand, S.B., C. S. Ravishakar & E. J. Delp. 1991. An iterative growing and pruning algorithm for classification tree design. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13, 167-174.
- [6] Hadiyanto, S. 1994. Metode Prakiraan. Badan Meteorologi Dan Geofisika Balai Wilayah II, Jakarta (tidak dipublikasikan).
- [7] Keperta, S. 1996. Non-binary classification trees. *Statistics and Computing*. 06:231-243.
- [8] Mc. Hattie, L.B. & F. Schnelle. 1974. *An introduction to Agrotopoclimatology*. WMO. TN. No. 133. Geneva.
- [9] Therneau, T.M. & E. J. Atkinson. 1997. An Introduction to Recursive Partitioning using the RPART routines. *Technical Report; Mayo Foundation*.
- [10] Venables, W. N. & B. D. Ripley. 1994. *Modern Applied Statistics with S-plus*. Springer-Verlag, New York. Inc.
- [11] Zhang, Heping. 1998. Classification trees for multiple binary responses. *Journal of the American Statistical Association*. 93: 180-193.

