

Grammatical Evolution Tatami untuk Ekstraksi Fitur dengan Pengukuran Multi Fitness

Go Frendi Gunawan¹⁾

Joko Lianto Buliali²⁾

1) Jurusan Informatika, Fakultas Teknologi Informasi, Kampus ITS Keputih, Sukolilo, Surabaya,
email: gofrendiasgard@gmail.com

2) Jurusan Informatika, Fakultas Teknologi Informasi, Kampus ITS Keputih, Sukolilo, Surabaya,
email: joko@its-sby.edu

Abstract: Feature Extraction is a significant topic in classification problem. Until now, there is no standard way to determine best features of data. In this thesis, grammatical evolution with multiple fitness evaluation approach (named as GE Tatami) has been developed to extract best features of data. The method generates $n-1$ features to separate data hierarchically, with n is number of classes. Some methods have been evaluated in this research, including genetics algorithm, grammatical evolution with global fitness measurement, grammatical evolution with multi fitness measurement, grammatical evolution with Tatami fitness measurement, and Gavrilis's grammatical evolution. It is shown in the experiment that Tatami method produces better results compared to the four other methods for synthesis data using decision tree classifier. The synthesis data is hierarchically separable. However Tatami method fails to boost SVM's accuracy. This method also fails when ideal features cannot be found..

Keywords: feature -extraction, grammatical evolution, classification, multi-fitness, Tatami.

Ekstraksi fitur adalah proses untuk menemukan pemetaan dari fitur-fitur asli ke dalam fitur-fitur baru yang diharapkan dapat menghasilkan keterpisahan kelas secara lebih baik (Guo, 2011). Ekstraksi fitur merupakan topik penting dalam klasifikasi, karena fitur-fitur yang baik akan sanggup meningkatkan tingkat akurasi, sementara fitur-fitur yang tidak baik cenderung memperburuk tingkat akurasi.

Beberapa metode berbasis algoritma genetika telah dibuat dalam penelitian-penelitian sebelumnya guna mencari fitur-fitur terbaik. Pada (Gavrilis, 2005) dan (Gavrilis, 2008), telah dikembangkan metode grammatical evolution untuk kepentingan ekstraksi fitur. Pada penelitian tersebut, Gavrilis dkk telah membuat metode yang berhasil membuat fitur-fitur baru dengan memanfaatkan akurasi classifier sebagai *fitness function*. Pada (Guo, 2011) dan (Li 2011), metode yang hampir sama juga digunakan untuk kasus yang berbeda. Adapun dalam penelitian-penelitian tersebut, terkadang fitur-fitur yang tidak relevan juga ikut tercipta dan digunakan sebagai dalam proses klasifikasi.

Pada (Gunawan, 2012), diciptakan pendekatan baru dengan mengukur akurasi setiap fitur guna meminimalisasi kemungkinan terlibatnya fitur-fitur yang tidak relevan dalam proses klasifikasi. Adapun pendekatan ini

cenderung menghasilkan akurasi yang buruk, karena dalam metode tersebut hanya dihasilkan 1 fitur untuk memisahkan semua kelas.

Dalam penelitian ini, diusulkan pendekatan multi fitness untuk memisahkan data secara hirarkikal. Pendekatan ini dinamai Grammatical Evolution Tatami (GE Tatami) karena dalam metode ini tampak bentuk yang menyerupai rantai tradisional Jepang (tatami). Sebagai pembanding, beberapa metode sebelumnya juga turut disertakan.

DATA DAN RUANG FITUR

Pada permasalahan klasifikasi, umumnya data terdiri dari sejumlah baris, kolom (fitur), dan kelas. Data numerik pada tabel 1 terdiri dari 2 fitur (x dan y) serta 3 kelas (A, B, dan C).

Karena data tersebut terdiri dari 2 fitur asli, maka dimungkinkan untuk menggambarannya dalam diagram cartesius untuk menghasilkan visualisasi yang lebih baik. Pada gambar 1 disajikan visualisasi dari data di tabel 1.

Seperti yang tampak pada gambar 1, mustahil untuk memisahkan data berdasarkan kelas hanya dengan menggunakan salah satu fitur asli saja (x atau y). Kedua fitur tersebut harus digunakan secara bersama-sama untuk

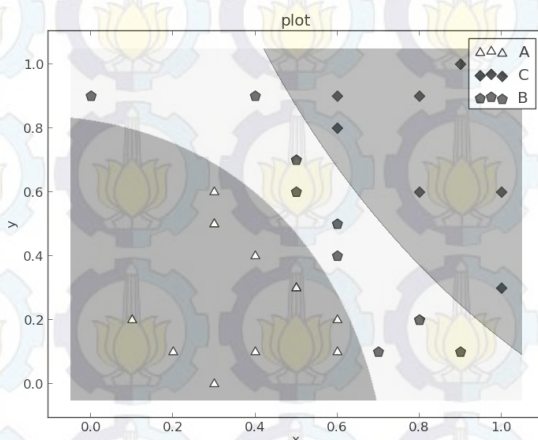
menghasilkan ruang fitur yang memungkinkan terjadinya pemisahan data.

Classifier yang umum digunakan semisal SVM atau jaringan syaraf tiruan akan dapat memisahkan data berdasarkan kelas dengan menggunakan ruang fitur yang disajikan dalam gambar 1 dengan cukup baik. Adapun demikian, pemisahan data tidak bisa dilakukan secara linear. Garis-garis pemisah yang dihasilkan akan berbentuk kurva yang secara matematis lebih rumit dibandingkan sekedar garis linear.

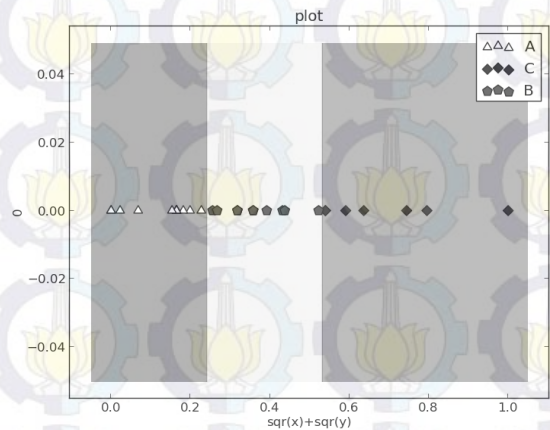
Pada permasalahan klasifikasi, adalah wajar untuk memetakan data yang tidak terpisahkan (atau sulit terpisahkan) ke dalam dimensi yang lebih tinggi. Dimensi yang lebih tinggi berkecenderungan untuk memberikan peluang keberhasilan yang lebih tinggi pula dalam memisahkan data. Adapun demikian, pemetaan ke dimensi yang lebih tinggi juga akan berimplikasi pada perhitungan yang lebih rumit.

Tabel 1. Contoh Data Numerik.

Original Features		Class
x	y	
0.3	0.5	A
0.4	0.9	B
0.6	0.2	A
0.9	1.0	C
1.0	0.3	C
0.8	0.2	B
...



Gambar 1. Ruang Fitur yang Dibentuk dari Fitur Asli dan Garis Pemisah yang Dibentuk oleh SVM dengan Kernel RBF



Gambar 2. Ruang Fitur yang Dibentuk oleh Fitur Baru $(\sqrt{x} + \sqrt{y})$

EKSTRAKSI FITUR

Pada dasarnya ekstraksi fitur adalah pemetaan data dari fitur-fitur asli ke bentuk lain. Proses pemetaan tersebut mungkin membentuk dimensi yang lebih tinggi atau rendah daripada ruang fitur asli.

Tujuan ekstraksi fitur adalah untuk membuat sesedikit fitur yang sanggup menciptakan pemisahan data dengan cara yang relatif sederhana (misalnya dengan menggunakan garis-garis linear).

Untuk lebih menjelaskan tujuan dan kegunaan dari ekstraksi fitur, data pada tabel 1 yang terdiri dari 2 fitur asli (x dan y) ditransformasi ke dalam sebuah ruang fitur 1 dimensi yang hanya berisi 1 fitur baru $(\sqrt{x} + \sqrt{y})$. Ruang fitur yang terbentuk disajikan pada gambar 2.

Pada hasil transformasi tersebut, dimungkinkan untuk memisahkan data dengan menggunakan dua buah garis linear yang secara matematika jauh lebih sederhana daripada sebelumnya (pada gambar 1).

GRAMMATICAL EVOLUTION

Grammatical Evolution adalah sebuah algoritma evolusi berbasis algoritma genetik. Metode ini memanfaatkan context-free *grammar* terdefinisi untuk mengubah suatu individu menjadi bentuk apapun yang dimungkinkan oleh *grammar* tersebut.

Dalam *grammatical evolution*, setiap individu terdiri dari string biner atau integer yang disebut genotip. Dengan memanfaatkan *grammar* terdefinisi, genotip tersebut diubah menjadi fenotip. Fenotip yang terbentuk dapat

berupa operasi matematika, sebuah fungsi, atau bahkan kode program komputer lengkap, tergantung dari *grammar* yang digunakan.

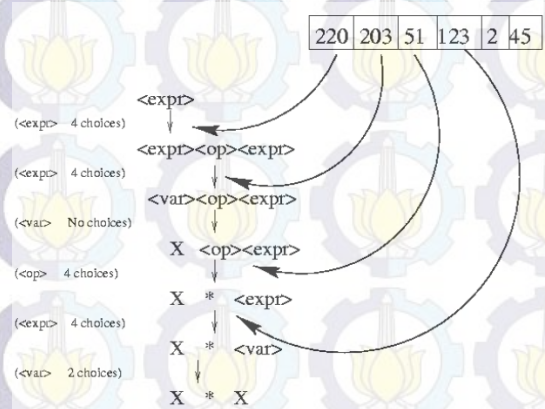
Sebagai contoh, didefinisikan *grammar* seperti pada tabel 2. Semisal genotip sebuah individu terdiri dari integer string 220, 203, 51, 123, 2, 45. Proses dimulai dari node $\langle expr \rangle$ sebagai start symbol. $\langle expr \rangle$ memiliki 4 kemungkinan aturan evolusi. Karena itu diambil segment pertama dari genotip (dalam kasus ini 220), dan dilakukan operasi modulo. Karena $220 \bmod 4$ menghasilkan nilai nol, maka digunakanlah aturan ke nol. Hal ini mengakibatkan node $\langle expr \rangle$ berevolusi menjadi $\langle expr \rangle \langle op \rangle \langle expr \rangle$ sesuai aturan ke nol. Selanjutnya diambil segment kedua dari genotip (dalam kasus ini 203) dan node non-terminal pertama dari calon fenotip $\langle expr \rangle \langle op \rangle \langle expr \rangle$ (dalam kasus ini $\langle expr \rangle$). Sekali lagi dilakukan operasi modulo. Karena $\langle expr \rangle$ memiliki 4 kemungkinan aturan evolusi, dan $203 \bmod 4$ menghasilkan angka 3, maka diambillah aturan ketiga (dalam hal ini $\langle var \rangle$). Sekarang calon fenotip berubah menjadi $\langle var \rangle \langle op \rangle \langle expr \rangle$. Proses dilanjutkan sampai diperoleh $x * x$ sebagai fenotip individu. Gambar 3 menunjukkan proses transformasi secara lengkap.

Setelah mendapatkan fenotip, proses dilanjutkan seperti halnya pada algoritma genetika. Nilai *fitness* fenotip diukur dengan menggunakan *fitness function* tertentu, menghasilkan *fitness value*. *Fitness value* dari setiap individu digunakan untuk menentukan survival rate dari setiap individu, yang akan berkorelasi dengan keterpilihan individu terkait dalam generasi selanjutnya.

Tabel 2. Contoh Grammar

Node	Production Rule	Index
$\langle expr \rangle$	$\langle expr \rangle \langle op \rangle \langle expr \rangle$	0
	$(\langle expr \rangle \langle op \rangle \langle expr \rangle)$	1
	$\langle pre-op \rangle \langle expr \rangle$	2
	$\langle var \rangle$	3
$\langle op \rangle$	+	0
	-	1
	*	2
	/	3
$\langle pre-op \rangle$	Sin	0
	Cos	1
	Tan	2
var	x	0

Sumber: Conor, 2006



Gambar 3. Proses Transformasi
Sumber: Conor, 2006.

Untuk penyelesaian masalah ekstraksi fitur, umumnya akurasi *classifier* digunakan sebagai *fitness value*. Oleh sebab itu *fitness value* dari setiap individu akan berkisar antara 0 sampai 1.

Seperti halnya dalam algoritma genetika, pada *grammatical evolution*, akan ditemukan individu terbaik setelah beberapa generasi. Individu terbaik ini kemudian diambil sebagai solusi optimal untuk masalah yang dipecahkan

GE Multi

Semisal sebuah data terdiri dari n kelas, adalah cukup logis untuk berpikir bahwa dapat digunakan n buah fitur untuk memisahkan data berdasarkan kelasnya. Masing-masing fitur akan bertugas untuk memisahkan satu kelas tertentu dengan semua kelas lainnya.

Untuk mencapai pendekatan ini, sebuah *grammatical evolution* sederhana dengan satu nilai *fitness* untuk setiap individu tidak akan dapat bekerja dengan baik. Oleh sebab itu, dibuatlah versi modifikasi dari *grammatical evolution* yang menggunakan sejumlah nilai *fitness* untuk setiap individu. Versi termodifikasi ini diberi nama GE Multi.

Semisal terdapat n kelas: $\{c_1, c_2, c_3, \dots, c_n\}$. Dalam GE Multi, setiap individu akan memiliki n *fitness value*: $\{f_1, f_2, f_3, \dots, f_n\}$. *Fitness value* pertama f_1 menggambarkan tingkat keberhasilan fenotip untuk memisahkan kelas c_1 dan semua kelas lainnya ($\{c_2, c_3, \dots, c_n\}$). *Fitness value* f_2 menggambarkan tingkat keberhasilan fenotip untuk memisahkan c_2 dan semua kelas lainnya ($\{c_1, c_3, \dots, c_n\}$). *Fitness value* ke n , f_n menggambarkan tingkat keberhasilan fenotip untuk memisahkan c_n dan semua kelas lainnya ($\{c_1, c_2, c_3, \dots, c_{n-1}\}$).

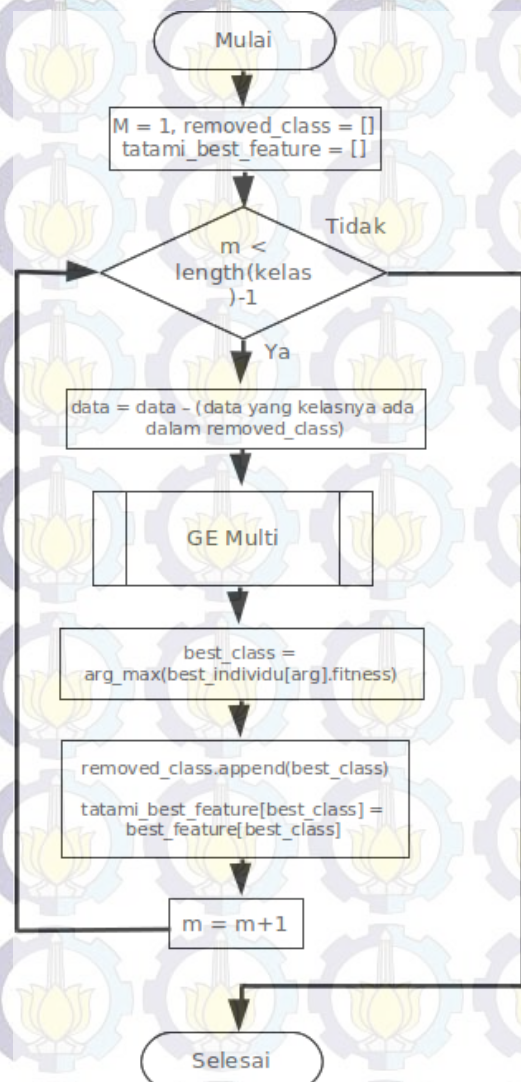
Untuk menghitung nilai *fitness* dari setiap individu, data harus ditransformasi menjadi 2 kelas. Untuk mengukur nilai *fitness* pertama f_1 , maka kelas-kelas $\{c_2, c_3, \dots, c_n\}$ harus digabungkan menjadi satu kelas (dimisalkan c_{-1}). Kemudian *classifier* akan bertugas untuk memisahkan c_1 and c_{-1} . Akurasi *classifier* kemudian digunakan sebagai *fitness value* f_1 .

Dengan menggunakan pendekatan ini, maka akan terbentuk n atau kurang individu-individu terbaik. Setiap individu memiliki nilai *fitness* terbaik untuk setiap kelas (akan ada individu dengan nilai f_1 terbaik, individu dengan nilai f_2 terbaik, dan seterusnya). Ada kemungkinan pula, bahwa sebuah individu memiliki beberapa nilai *fitness* terbaik. Dengan demikian, penggunaan GE Multi sebagai *feature extractor* akan menghasilkan maksimum n fitur baru dimana n adalah jumlah kelas dalam data.

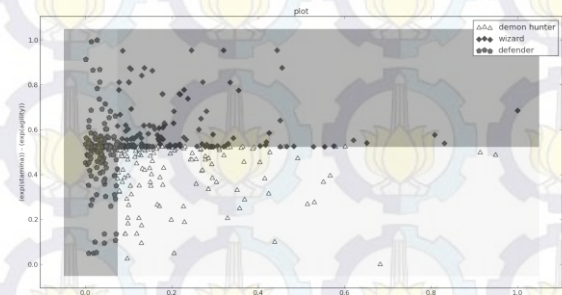
GE Tatami

GE Tatami merupakan pengembangan dari GE Multi yang didesain untuk data yang terpisah secara hirarkikal. Untuk data yang terpisah secara hirarkikal, maka pemisahannya pun akan dapat dilakukan secara hirarkikal. Sebagai contoh, semisal ada n kelas, $\{c_1, c_2, c_3, \dots, c_n\}$. Sekali c_1 terpisah dari kelas-kelas lainnya, maka dalam proses selanjutnya c_1 dapat diabaikan, sehingga proses pemisahan hanya akan terfokus pada $\{c_2, c_3, \dots, c_n\}$.

GE Tatami terdiri dari $n-1$ iterasi di mana n =jumlah kelas. Pada setiap iterasi, GE Multi digunakan untuk menentukan individu-individu terbaik. Nilai *fitness* maksimum dari semua nilai *fitness* terbaik dari GE Multi digunakan untuk menentukan kelas mana ya paling mungkin terpisah dari semua kelas la. Individu yang terkait dengan kelas yang paling terpisah kemudian dipilih dan dimasukkan dalam list fitur baru. Kelas yang sudah terpisah itu selanjutnya dihilangkan dari data. Proses iterasi selanjutnya dilanjutkan dengan kelas ya telah dihilangkan. Dalam setiap iterasi, jumlah data akan terus berkurang sesuai dengan kelas yang telah terpisah. Pada iterasi terakhir hanya akan tersisa 2 kelas. Proses berakhir setelah ditemukan individu terbaik untuk memisahkan kedua kelas terakhir tersebut. Gambar 4 menunjukkan flowchart global dari GE Tatami.

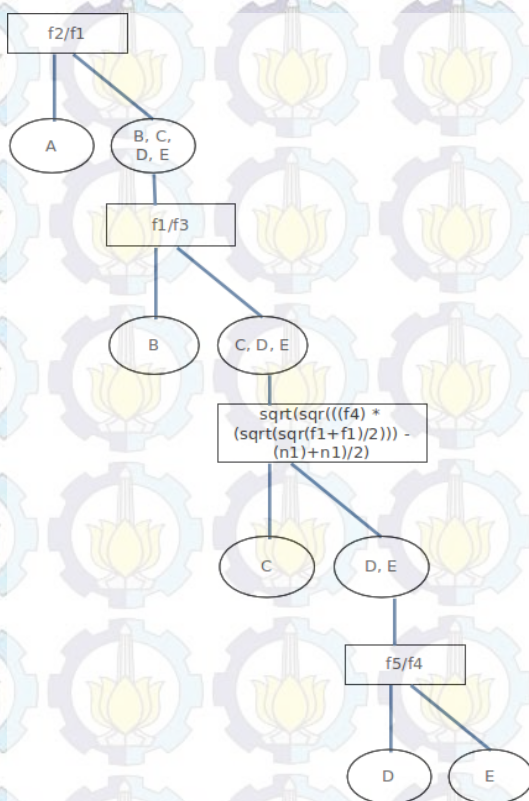


Gambar 4. Flowchart Global GE Tatami



Gambar 5. Ruang Fitur yang Dihasilkan GE Tatami

Dengan pendekatan ini, maksimum akan didapatkan $n-1$ fitur baru. Gambar 5 menunjukkan bagaimana dua buah fitur yang dihasilkan oleh GE Tatami sanggup memisahkan tiga kelas dengan sangat baik.



Gambar 6. Decision Tree yang Dihasilkan dengan Menggunakan Fitur-Fitur yang Dihasilkan GE Tatami

GE Tatami sangat cocok sebagai langkah pre-processing pada *classifier decision tree*, karena keduanya sama-sama dirancang untuk memisahkan data secara hirarkikal. *Decision Tree* yang terbentuk dengan menggunakan fitur-fitur yang dihasilkan oleh GE Tatami ditunjukkan pada gambar 6.

HASIL DAN PEMBAHASAN

Untuk membandingkan efektifitas GE Multi dan GE Tatami dengan metode-metode sebelumnya, telah dibuat beberapa percobaan yang melibatkan sejumlah dataset. Dataset yang digunakan dalam percobaan terdiri dari 3 buah dataset sintesis yang terpisah secara hirarkikal serta 3 dataset lain yang umum diambil dari website UCI Machine Learning (iris, balanced scale, dan E-Coli).

Pada dataset sintesis 03, fitur pemisah sengaja disembunyikan dalam beberapa fitur tampak. Ini bertujuan untuk menguji kemampuan GE Tatami dalam memisahkan data yang secara eksplisit tidak tampak terpisah secara hirarkikal, namun sebenarnya terdapat cara untuk memisahkan data tersebut secara

hirarkikal. Fitur-fitur pemisah dapat diperoleh kembali dengan melakukan otak-atik menggunakan operasi matematika pada fitur-fitur tampak.

Untuk setiap dataset, dilakukan 5 fold cross-validation dan sebuah test lain yang melibatkan semua data sebagai training sekaligus testing. Hasil lengkap percobaan dapat diakses secara bebas pada alamat <https://github.com/goFrendiAsgard/feature-extractor/>

Hasil percobaan dirangkum dengan menggunakan perhitungan rata-rata untuk setiap dataset, menghasilkan tabel 3 (keterangan: A=akurasi dalam persen, F=Jumlah Fitur, Trn=Training, Tst=Testing, Ttl=Total). Dengan melihat pada tabel tersebut, tampak bahwa GE Tatami menunjukkan hasil yang sangat baik pada data sintesis 02 dan sintesis 03, serta hasil yang cukup baik pada dataset sintesis 01. Adapun demikian, GE Multi menunjukkan akurasi yang lebih baik pada data-data non-sintesis.

Tabel 3 Rata-Rata Hasil Percobaan

Percobaan	GA Select Feature	GE Global		GE Multi		GE Tatami		GE Gavrilis	
		A	F	A	F	A	F	A	F
Iris	Trn	96	2	99	1	99	3	99	2
	Tst	94		93		95		93	
	Ttl	96		98		98		98	
B-Scal	Trn	71	3	89	1	96	1	93	2
	Tst	69		79		82		82	
	Ttl	71		87		93		91	
E.Coli	Trn	97	5	87	1	98	8	97	7
	Tst	78		65		64		51	
	Ttl	93		83		92		88	
Syn-01	Trn	74	3	79	1	100	3	100	2
	Tst	72		42		85		82	
	Ttl	74		72		97		97	
Syn-02	Trn	78	4	71	1	100	3	100	3
	Tst	74		40		73		79	
	Ttl	78		65		95		96	
Syn-03	Trn	73	5	68	1	99	5	100	4
	Tst	41		37		63		78	
	Ttl	66		62		92		96	
AVG	Trn	82	3	82	1	99	3	98	3
	Tst	71		59		77		77	
	Ttl	80		78		94		94	

Keunggulan mutlak GE Tatami pada dataset sintesis 02 dan sintesis 03 terjadi karena dataset tersebut terpisah secara hirarkikal, dan GE Tatami berhasil menemukan ruang fitur ideal untuk *classifier decision tree* (seperti yang digambarkan dalam gambar 5 dan gambar 6). Di sisi lain, GE Tatami tidak menunjukkan keunggulan mutlak untuk dataset sintesis 01. Walaupun dataset sintesis 01 juga terpisah secara hirarkikal, namun data tersebut hanya memiliki 3 kelas sehingga tidak memiliki struktur hirarkikal yang cukup dalam, dan karenanya menjadi kurang signifikan dalam proses.

SIMPULAN

Hasil percobaan menunjukkan bahwa GE Tatami menunjukkan performa yang baik pada data yang terpisah secara hirarkikal dengan menggunakan *classifier decision tree*, khususnya data yang terdiri dari banyak kelas (seperti data sintesis 02 dan sintesis 03). Adapun demikian, menggunakan *classifier* yang sama, GE Multi menunjukkan hasil yang relatif baik dan seimbang untuk semua dataset.

Dari hasil tersebut, dimunculkan hipotesa bahwa menggabungkan fitur-fitur yang dibentuk oleh GE Multi dan GE Tatami mungkin akan menghasilkan akurasi yang baik untuk data dengan cakupan karakteristik yang lebih luas (tidak hanya yang terpisah secara hirarkikal).

Dalam percobaan yang dilakukan, akurasi *classifier* digunakan untuk menentukan nilai *fitness* dari setiap individu. Penggunaan mekanisme yang lebih baik dan sederhana dalam pengukuran *fitness value* mungkin akan mempercepat proses ekstraksi fitur

UCAPAN TERIMA KASIH

Beberapa orang telah berjasa dalam membantu penulis utama untuk menyelesaikan penelitian ini. Salah satu yang paling berperan penting adalah Bapak Joko Lianto, sebagai pembimbing dalam penelitian ini. Berbagai saran yang beliau berikan telah menginspirasi penulis utama untuk menyelesaikan berbagai macam permasalahan dan menyelesaikan penelitian ini tepat waktu.

Beberapa teman, seperti Bapak Mukhlis Amien dan Bapak Hendra Suprayogi, juga turut berperan dalam memberikan saran-saran yang bersifat *out-of-the-box* dan terkadang menjadi sangat berguna dalam penyelesaian masalah.

RUJUKAN

- Conor, R. (2006). Grammatical Evolution Tutorial. Gecco 2006.
- Gavrilis D., Tsoulous I. G., Georgoulas G., Glavas E. (2005). "Classification of Fetal Heart Rate Using Grammatical Evolution", IEEE Workshop on Signal Processing Systems Design and Implementation, 2005.
- Gavrilis D., Tsoulous I. G., Dermatas E. (2008). "Selecting and Constructing Features Using Grammatical Evolution", Journal Pattern Recognition Letters Volume 29 Issue 9, July, 2008 Pages 1358-1365.
- Guo L., Rivero D., Dorado J., Munteanu C. R., Pazos A. (2011). "Automatic feature extraction using genetic programming: An application to epileptic EEG classification", Expert Systems with Applications 38 Pages 10425-10436
- Gunawan G. F., Gosaria S., Arifin A. Z. (2012). "Grammatical Evolution For Feature Extraction In Local Thresholding Problem", Jurnal Ilmu Komputer dan Informasi, Vol 5, No 2 (2012)
- Harper R., Blair A. (2006). "Dynamically Define Functions in Grammatical Evolution", IEEE Congress of Evolutionary Computation, July 16-21, 2006
- Li B., Zhang P.Y., Tian H., Mi S.S., Liu D.S., Ruo G.Q. (2011). "A new feature extraction and selection scheme for hybrid fault diagnosis of gearbox", Expert Systems with Applications 38 Pages 10000-10009
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, E. (2011). "Scikit-learn: Machine Learning in Python", Journal of Machine Learning Research Vol. 12 Pages 2825-2830 Trans. Roy. Soc. London, vol. A247, pp. 529-551, April 1955.