

**PENERAPAN TEKNIK SVM (*Support Vector Machine*) UNTUK
PENDETEKSIAN PENGARUH KESTABILITASAN ENZIM
TERMUTASI**

KARYA AKHIR

**Diajukan sebagai salah satu syarat untuk memperoleh gelar Magister
Teknologi Informasi**

SUSANTI KUSUMAWIDJAYA

0706193920



**UNIVERSITAS INDONESIA
FAKULTAS ILMU KOMPUTER
PROGRAM STUDI MAGISTER TEKNOLOGI INFORMASI**

DEPOK

JANUARI 2009

ABSTRACT

A basic question in protein science is to which extent mutations affect protein thermostability or not. This knowledge would be particularly relevant for engineering thermostable enzymes. In several experimental approaches, this issue has been serendipitously addressed. The big problem is, this manual experimental will waste high budget because it needs trial and error manual experimental procedure. Trial and error procedural on the laboratory are used to find the best fits procedural that will use to make the experimental succes, but needs more time and more budget. It would be therefore convenient providing a computational method that predicts when a given protein mutant is more thermostable than its corresponding wild-type. The purpose of this research is to compare the accuracy data from SVM with the data from experimental manual that was done by Annaluru(2002) and Salazar(2007). SVM (Support Vector Machine) is a machine learning methode that will used in this research as a predictor to predict the enzyme thermostability. BLAST (Basic Local Alignment Tools) are used as a data mining methode to extract data from 14 family protein (Annaluru and Salazar data experimental). From the outputs of the BLAST runs only aligned sequence pairs comprising a mesophilic and a thermophilic protein were selected. Only protein pairs sharing at least 70% of sequence identity were retained (a redundancy reduced dataset). we used SVM L20 (difference of the residue composition in each pair of the dataset) and SVM L400 (difference of the dipeptide composition in each pair of the dataset). When trained and tested on a redundancy reduced dataset (homology >70%), our predictor (SVM) achieves 86% accuracy for SVM L20 (classifies 12 out of 14 experimentally characterized protein mutants with enhanced thermostability from Annaluru and Salazar) and 79% for SVM L400 (classifies 11 out of 14 experimentally characterized protein mutants with enhanced thermostability from Annaluru and Salazar). This accuracy from SVM are one of the reasonable way to used SVM as a methode to reduce the trial and error procedural because can give the prediction and reduce the failed experimental.

ABSTRAK

Salah satu hal mendasar di dunia science mengenai enzim adalah, penelitian untuk memprediksi pengaruh perpanjangan mutasi terhadap thermostabilitas enzim yang seringkali berkaitan dengan penelitian rekayasa enzyme thermostabil. Masalah utama muncul ketika landasan untuk melakukan percobaan secara manual di laboratorium, prosedural pengerjaannya selalu dilakukan secara trial and error. Dana dan waktu yang terbuang untuk mendukung percobaan ini cukup besar, mengingat percobaan masih dilakukan secara manual dan juga faktor keberhasilan yang belum akurat. Hal ini yang mendasari diperlukannya suatu campur tangan secara komputerisasi untuk memberikan prediksi mengenai thermostabil enzim mutant yang diperbandingkan dengan *wild type*-nya(sebelum dimutasi). Tujuan dari penelitian ini adalah untuk memperbandingkan keakurasian data yang diperoleh dari SVM terhadap data experimental yang dilakukan secara manual oleh Annaluru(2002) dan Salazar(2007). SVM (*Support Vector Machine*) adalah suatu metode *Machine Learning* yang digunakan dalam penelitian ini untuk memberikan prediksi stabilitas enzyme termutasi. Data set yang digunakan untuk SVM berasal dari data yang di hasilkan sebagai output dari BLAST (*Basic Local Alignment Tools*). BLAST digunakan sebagai metode untuk me-mining data dari ke-14 family protein (Annaluru dan Salazar experimental data). Mining data dari ke-10 protein ditujukan untuk mendapatkan kumpulan protein mutant yang memiliki tingkat homologi >70% dari sekumpulan besar data 10 family protein yang belum jelas tingkat homologinya. Proses SVM dilakukan dua kali, yaitu SVM L20 Untuk residu asam amino yang berbeda dalam setiap pasang data *sequence* dalam data set dan SVM L400 Untuk *sequence* dengan komposisi dipeptida yang berbeda. Hasil akhir, SVM memberikan tingkat keakurasian yang tidak jauh berbeda, yaitu sebesar 86% untuk L20 (SVM berhasil mengidentifikasi 12 dari 14 jenis protein) dan 79% untuk L400 (SVM berhasil mengidentifikasi 11 dari 14 jenis protein). keakurasian SVM yang hampir memprediksi dengan sempurna mengenai thermostabilitas enzyme termutasi, dapat digunakan sebagai metode untuk meminimalisir kesalahan dan waktu yang terbuang dalam trial dan error prosedural experimental.

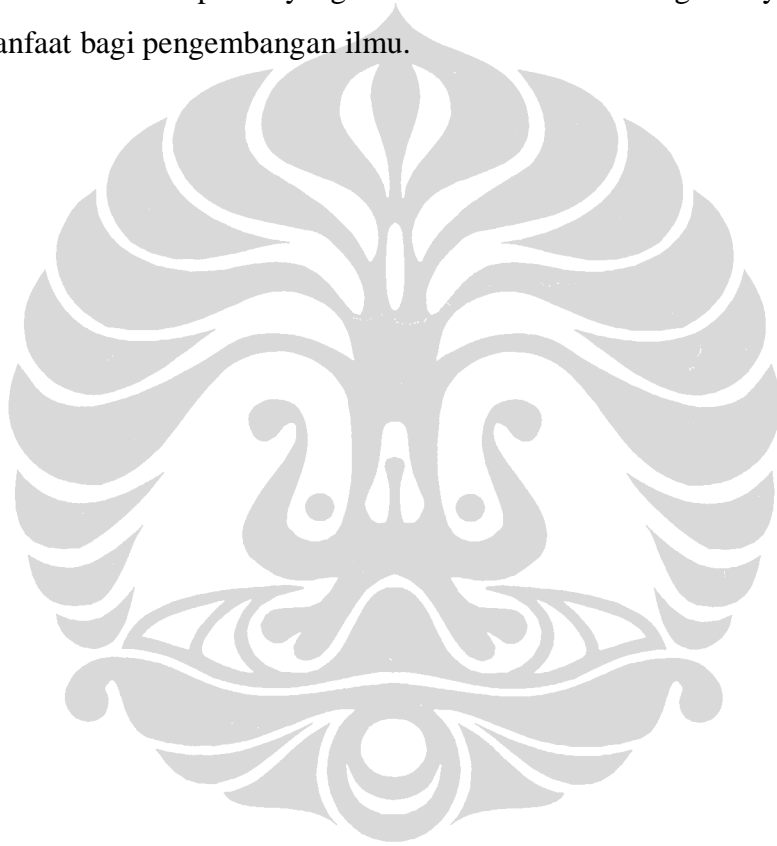
KATA PENGANTAR

Puji syukur penulis panjatkan ke hadirat Allah SWT atas segala rahmat dan karunia-Nya yang senantiasa dicurahkan kepada umat-Nya khususnya penulis, sehingga penulis dapat menyelesaikan skripsi ini. Shalawat serta salam semoga selalu tercurah kepada Rasulullah SAW selaku suri tauladan yang baik dan kepada para sahabat, keluarga dan pengikutnya hingga akhir zaman. Penulisan karya akhir ini dilakukan dalam rangka memenuhi salah satu syarat untuk mencapai gelar Magister Teknologi Informasi pada Program Studi Magister Teknologi Informasi, Fakultas Ilmu Komputer - Universitas Indonesia. Saya menyadari bahwa, tanpa bantuan dan bimbingan dari berbagai pihak, dari masa perkuliahan sampai pada penyusunan karya akhir ini, sangatlah sulit bagi saya untuk menyelesaikannya. Oleh karena itu, saya mengucapkan terima kasih kepada:

- (1) Prof Dr Aniati Murni selaku pembimbing yang dengan sabar meluangkan waktunya untuk memotivasi dan memberikan petunjuk mengenai apa yang penulis perlukan untuk menyelesaikan thesis ini.
- (2) Prof Dr Soetijoso Soemitro selaku responden yang turut membantu untuk memberikan masukan kepada penulis mengenai materi bioinformatika dan protein.
- (3) Para Staff Laboratorium Biokimia LIPA UNPAD yang telah banyak membantu dalam usaha memperoleh data yang saya perlukan.
- (4) Kedua orang tua penulis yang telah memberikan dukungan moril dan dukungan materil yang luar biasa.
- (5) Seluruh dosen dan karyawan MTI-UI atas semua yang telah diberikan selama penulis menjalani masa perkuliahan.

- (6) Teman-teman satu kelas angkatan 2007SA atas kekompakannya, kerja sama dan kebersamaan yang dibangun selama kuliah.
- (7) Serta semua pihak yang telah membantu terlaksananya pembuatan thesis ini.

Akhir kata, saya berharap Tuhan Yang Maha Esa berkenan membalas segala kebaikan semua pihak yang telah membantu. Semoga karya akhir ini membawa manfaat bagi pengembangan ilmu.



Jakarta, Januari 2009

Penulis

DAFTAR ISI

ABSTRAK.....	i
KATA PENGANTAR.....	ii
DAFTAR ISI.....	iv
DAFTAR TABEL.....	vi
DAFTAR GAMBAR.....	vii
BAB I PENDAHULUAN.....	1
1.1 LATAR BELAKANG.....	1
1.2 PERUMUSAN MASALAH.....	4
1.3 TUJUAN DAN MANFAAT PENELITIAN.....	5
1.4 RUANG LINGKUP.....	6
1.5 SISTEMATIKA PENULISAN.....	8
BAB II DATA MINING DALAM BIOINFORMATIKA.....	9
2.1 DEFINISI DATA MINING.....	9
2.2 BLAST SEBAGAI TOOLS DATA MINING UNTUK BIOINFORMATIKA.....	12
2.1.1 NCBI Sebagai Pusat Data Untuk Sequence Protein.....	13
2.2.2 BLAST Sebagai Metode Untuk me-Mining Data Sequence Protein.....	17
2.2.1.1 BLAST <i>all-against-all</i>	8
2.2.1.2 Kesulitan dan kendala yang mempengaruhi tingkat “homologi” BLAST pada urutan asam amino.....	21
2.3 SVM (<i>Support Vector Machine</i>).....	23
2.3.1 <i>Structural Risk Minimization (SRM)</i>	25
2.3.2 SVM pada <i>Linearly Separable Data</i>	25
2.3.3 SVM pada <i>Nonlinearly Separable Data</i>	26

2.3.4	Multi Class SVM.....	27
2.3.5	One Class SVM.....	28
2.3.6	<i>Estimasi Parameter</i>	29
2.3.7	Supervised Learning.....	30
2.4	SVM UNTUK MEMPREDIKSI THERMOSTABILITAS ENZIME TERMUTASI.....	31
BAB III PENERAPAN TI PADA <i>LABORATORIUM RESEARCH ENZIM</i>		
	THERMOSTABIL.....	35
3.1	ENZIM THERMOSTABIL DAN INDUSTRI.....	37
	3.1.1 Studi penyebab Enzime Thermofilik.....	38
	3.1.2 Keuntungan Enzime Thermostabil.....	39
3.2	Produksi Enzime di Indonesia.....	40
3.3	MUTASI.....	42
	3.3.1 Penelitian mengenai Mutasi.....	42
	3.3.2 Frekuensi mutasi terhadap jenis asam amino.....	44
3.4	PERMASALAHAN MANUAL EXPERIMENT.....	45
3.5	PERANAN TI DALAM BIOINFORMATIKA.....	50
BAB IV METODOLOGI PENELITIAN.....		
4.1	LANGKAH-LANGKAH PENELITIAN.....	55
	4.1.1 Keterangan Tahapan Kerja.....	59
	4.2.1 <i>Random Data</i>	67
4.2	DATA DAN PERALATAN PENELITIAN.....	68
BAB V PEMBAHASAN HASIL.....		
5.1	DATA.....	70
	5.1.1 Data Preparation.....	71
5.2	PENGUKURAN EVALUASI PARAMETER.....	72

5.3	TRAINING DATA SET.....	73
5.4	<i>SCORING THE SVM PREDICTOR</i>.....	74
5.5	HASIL PENGUJIAN SVM.....	77
5.6	VALIDITAS DATA PENGUJIAN SVM DENGAN PENGUJIAN ANNALURU dan SALAZAR	79
	5.6.1 Penerapan IT dalam Laboratorium Experimental.....	80
BAB VI KESIMPULAN.....		83
6.1	KESIMPULAN.....	83
6.2	SARAN.....	84
DAFTAR PUSTAKA.....		85
LAMPIRAN A		
LAMPIRANB		



DAFTAR TABEL

Tabel 2.1	Mutasi asam amino E (Glu) pada posisi 200 dengan 19 asam amino dengan menggunakan <i>I-Mutant2.0</i>	32
Tabel 2.2	<i>Data set used to train and test the SVM exploiting structural information</i>	33
Tabel 3.1	Jenis-jenis penggunaan enzyme untuk industri.....	36
Tabel 3.2	<i>Frequency of Amino Acid Residues in Mutation Sites</i>	44
Tabel 4.1	Data 10 family Protein yang di mutasi oleh Annaluru dan Salazar.....	51
Tabel 4.2	Hasil pengujian Thermostabilitas 14 jenis protein oleh Annaluru dan Salazar.....	52
Table 5.1	<i>Performances obtained with random splitting of the cross-validation sets</i>	74
Table 5.2	<i>Performances of the two SVMs Method</i>	74
Tabel 5.3	<i>SVM performances for the experimental dataset</i>	76
Tabel 5.4	Perbandingan Hasil.....	77
Tabel 5.5	Hasil pengujian modul Pengenalan Karakter.....	78

DAFTAR GAMBAR

Gambar 2.1	Menu utama pada <i>Search Database</i> NCBI.....	14
Gambar 2.2	Bentuk Data pada NCBI untuk 1 protein.....	15
Gambar 2.3	Contoh urutan <i>sequence</i> asam amino dari NCBI.....	16
Gambar 2.4	Blast <i>sequence</i> asam amino 2 protein dalam 1 family.....	18
Gambar 2.5	Layout <i>BLAST all-against-all</i> strategy.....	19
Gambar 2.6	<i>Lay out</i> proses BLAST.....	20
Gambar 2.7	Hasil score dari <i>alignment Sfamy</i> dengan TVA1.....	22
Gambar 2.8	Perbandingan pelipatan (<i>folding</i>) domain C ALP1 dan SBD TVA1.....	23
Gambar 2.9	Alternatif bidang pemisah (kiri) dan <i>bidang pemisah</i> terbaik dengan margin (m) terbesar (kanan).....	29
Gambar 2.10	Contoh transformasi untuk data yang tidak dapat dipisahkan secara linier.....	26
Gambar 2.11	Transformasi ke <i>feature space</i>	29
Gambar 3.1	Tahapan <i>Enzyme Production</i>	41
Gambar 3.2	Perbandingan urutan asam amino domain C dari ALP1.....	43
Gambar 3.3	<i>Fish-bone</i> diagram.....	45
Gambar 3.4	Skema Alur Kerja Laboratorium Project Experiment.....	49
Gambar 4.1	Bagan alir Tahapan Penelitian.....	55
Gambar 4.2	Design Methode.....	57
Gambar 4.3	Ke-10 Family protein (acuan dari Annaluru dan Salazar).....	61
Gambar 4.4	Family dari protein β -GUS. Sumber diambil dari NCBI.....	62
Gambar 4.5	Potongan scoring Homology β -GUS.....	63
Gambar 5.1	Alur penentuan strategi experimental dengan dukungan IT.....	81