

MICHAEL PAUL STEWART BROWN

<http://www.cse.ucsc.edu/~mpbrown/>

430 Reina Del Mar Ave.
Pacifica, CA 94044
415-730-6211
michaelbrown14142@gmail.com

Summary

My goal is to participate in a high impact project concentrating on important real-world problems whose success relies upon the use of scientifically-rigorous statistical analysis, pattern recognition, and mathematical modeling techniques implemented using efficient computer algorithms and flexible data handling techniques. My ideal situation would be to work closely with a highly talented group of people.

I have been involved in a variety data analysis / pattern discovery projects, mostly concentrated on biological data and working with lab scientists to make biologically relevant discoveries: using proteomic mass spectrometry data to predict complex disease states like cancer, predicting new gene products in *Entamoeba histolytica*, using natural language in biomedical literature abstracts to predict the occurrence of names of organisms, using genomic sequence data to predict the best DNA markers to identify biowarfare pathogens, using video stream data to identify the presence of objects like people and vehicles, using customer data to rank customers to receive mail advertisements, using credit card data to predict fraudulent charges, using gene expression data to predict the functional classes of genes, using RNA sequences to predict the secondary structure of the molecule, ...

I have experience in many areas that help to build successful pattern recognition and data analysis projects. I have done theoretical development of models and algorithms including original work with hidden Markov models (HMMs), stochastic context-free grammars, and support vector machines to the application of biological data such as DNA sequences, microarray data, and proteomic mass spectrometry data. I have implemented theoretical ideas using efficient algorithms on various computer platforms including GNU/Linux compute clusters. I have performed large-scale experiments including the design, statistical analysis, and interpretation of them. I have experience collaborating with groups with diverse skill sets (computer science, biology, IT, business). I have experience in project management (timelines, milestones, deliverables, budgets), technology transfer, and business planning.

I enjoy pattern recognition and data analysis and am very fortunate to live in an era in which information is electronically available, the volume of information is growing exponentially (especially for biological data), and cheap computing power allows the automated analysis of it. This ultimately yields knowledge that enriches our lives.

Keywords: *statistical pattern recognition, Bayesian statistical analysis, statistics, support vector machines (SVMs), information retrieval, natural language processing, bioinformatics, genomics, gene expression analysis, proteomics, Hidden Markov Models (HMMs), stochastic context-free grammars, algorithm development, programming, Linux compute clusters, C++, C, Java, Perl, R, matlab, project management*

Education

- 1999** Ph.D Computer Science.
University of California, Santa Cruz (UCSC).
Title: “RNA Modeling Using Stochastic Context-Free Grammars”.
Advisor: David Haussler.

At UCSC, I was able to be a part of a spectrum of new work: hidden Markov models (HMMs) applied to sequence comparison, Dirichlet mixtures applied to protein priors, stochastic context-free grammars (SCFGs) applied to RNA secondary structure, support vector machines (SVMs) applied to DNA microarray gene expression data, SVMs on SCFG Fisher score vectors applied to organism identification, and others. All this work was done under a rigorous statistical pattern recognition and machine learning framework.

- 1992** B.S. Computer Science.
North Carolina State University (NCSU): Raleigh, NC.
Summa Cum Laude.

Past Positions

- **2006**
-

Stanford University Upi Singh Lab, Independent Collaborator

My work at the Singh lab has concentrated on the genomic analysis of the medically important parasite *Entamoeba histolytica*. This protozoan parasite causes colonic and liver disease that results in 100,000 deaths per year making it the second most common cause of parasitic death in humans. My genomic analysis has concentrated on identifying and characterizing spliced gene transcripts using large-scale alignments of EST libraries against the recently sequenced genome. I developed code to compute these alignments, store the large number of results in a database, identify possible spliced introns, correlate our findings with existing annotations, and explore new hypotheses regarding spliced gene products, alternatively spliced genes, and transposable elements. I also used hidden Markov models and stochastic context-free grammars to identify spliceosomal complex components in the full genome. I also maintain a webpage that details all the methods used, results and summaries of experiments, and allows interaction with the results database in a way that is easy to use for the biologists in the lab. This webpage has been crucial in making the collaboration with the biologists at Stanford possible as it is the easiest way to make available the large quantity of results that have been generated. The biological expertise of the lab coupled with my bioinformatic analysis has been fruitful. Several scientifically interesting findings are being documented in a paper that is in preparation.

- **2002-2005**
-

Predicant Biosciences, Sr. Scientist

My work at Predicant Biosciences involved the diagnosis of cancer disease states using mass spectrometry data collected from blood samples. My primary responsibility was to discover a pattern in protein abundance data that would differentiate healthy from cancer patients. Pattern discovery involved theoretical work as well as algorithm development. I worked daily with highly talented lab scientists to solve

critical problems.

The important challenge here is to find significant patterns that will hold up in an independent validation test set. This is challenging because the number of patients in the training set is relatively small (order 100) while the dimension of the raw data collected from each patient is large (order 10^6 , 10^8). The probability of overfitting is very significant. I used rigorous and well-founded techniques from statistical pattern recognition and machine learning. These techniques afford superior performance and an understanding of the data that any ad-hoc analysis technique cannot match. This is especially important in areas of medicine where an incorrect or misapplied analysis can have catastrophic consequences such as misdiagnosing a disease state like cancer.

I developed several pattern recognition techniques and implementations, feature selection wrappers, statistical tests, Quality Control tools, visualization software, and an automated pattern recognition system. Pattern recognition runs involving thousands of parameter sweeps were performed. Given a run specification, data servers were started, thousands of pattern recognition runs were made, the results were collected, and statistics were generated. This was done automatically and any result could be readily brought up for inspection.

In addition to the pattern discovery, I was intimately involved with the shear data processing burden of applying very complex algorithms to datasets whose total size could reach 60 gigabytes of data per day. This involved the use of a Linux compute cluster that ran cluster job control software that I developed. The analysis pipeline was complex that had about 6 steps with dependencies and would take about 4 hours of computation for each experiment. All this computation happened automatically as soon as the raw data was available and results were stored and available to everyone. QC decisions and system performance was gauged every day based on these computations.

Predicant was initially named Biospect.

• 1999-2002

Fair Isaac Corporation (FIC), Sr. Staff Scientist

I did most of my work in the Advanced Technologies Division of HNC. HNC was acquired by FIC in 2002.

I worked on several projects including: structural RNA identification working directly with lab scientists at ISIS Pharmaceuticals (funded through DARPA), a probabilistic model for information extraction from video streams (funded through ARDA), the bugID system for identifying optimal DNA probes for bacterial pathogens using a spectrum of data sources including biomedical literature abstracts and genomic sequences (funded through DARPA *principal investigator*), and identifying optimal DNA probes for viral pathogens (funded through DARPA and USAF *principal investigator*)

I was also involved with several projects involving natural language processing, information retrieval, fraud detection, and advertising optimization.

• 1996

Affymetrix, Researcher.

At Affymetrix, I worked on sequencing by hybridization using 10-mer DNA gene chips. I worked directly

with lab scientists whose expertise in chip chemistries and characterists helped solve many problems. The Affymetrix chip had every possible 10-mer on it and the problem was to do DNA *de novo* sequencing by using information about which probes were present in the sample.

Grants

2003 This grant was well reviewed but SBIR funding rules prevented us from accepting it. National Cancer Institute R44 CA105545-01 “Advanced Biological Pattern Discovery System”. Principal investigator.

We propose a discovery tool that accurately detects, diagnoses, and predicts outcomes of cancer using molecular signatures that characterize tumors and their interaction with their microenvironments and the body. One of the key components is a data analysis system that is robust and can reliably discover significant patterns from the large datasets produced by the separations-mass spectrometry components of our platform. Such a system will enable the discovery and assay of proteomic patterns that detect and distinguish different forms of cancer, including those that have similar clinical presentations, but differ at the molecular level. Predicant’s integrated platform and advanced informatics system will become vital for the informed clinical management of cancer.

This grant received an excellent priority score of 153. Scores range from 100 to 500 with 100 being the best. An excellent score usually means that it will be funded. However, due to the new interpretation of what constitutes a small business and the fact that the company was not owned by more than 51% individuals (being VC funded), we did not qualify to receive the SBIR grant money. The grant was scientifically sound but funding rules kept it from being funded.

2001 DARPA N66001-01-C-8010. “The BugID System for Discovering Optimal Nucleotide Probes”. Principal investigator.

The bugID system is a probabilistic method for discovering optimal DNA probes for pathogens. The system uses a wide variety of relevant information: genomic sequences from the pathogen, its close neighbors, and all other known organisms, the physical properties of the target molecule, phylogenetic relationships between organisms, biological truths mined from English text in biomedical literature abstracts, and physical and chemical properties of DNA hybridization. It combines all this information in a probabilistic framework to find optimal probes that have very high probability of hitting the target and very low probability of hitting any other organism sequence. Probes designed using this system have much better performance with fewer false positives and false negatives than a more simple system that does not account for all the relevant information. The software computes and keeps track of the large amount of information generated from the variety of sources, a complicated dataflow management problem.

1999 Program in Mathematics and Molecular Biology (PMMB) fellowship. Dr. D. Haussler.

1991 NSF Research Experience for Undergraduates. “Natural Language Processing in Scripted Domains”. Dr. R. Rodman.

Publications

1. Davis Carrie, Brown M.P.S., Singh U. (In preparation). *Entamoeba histolytica* Splicesomal Introns and snRNAs. *In preparation*
2. de Valpine P, Bitter HM, Brown MPS, Heller J. (submitted) Do biomarker discovery studies include enough patient samples to estimate valid patterns? *Submitted for publication*.
3. Sassi AP, Andel F 3rd, Bitter HM, Brown MP, Chapman RG, Espiritu J, Greenquist AC, Guyon I, Horchi-Alegre M, Stults KL, Wainright A, Heller JC, Stults JT. (2005). An automated, sheathless capillary electrophoresis-mass spectrometry platform for discovery of biomarkers in human serum. *Electrophoresis*, 26(7-8):1500-12.
4. Guyon I., Bitter H.M., Ahmed Z., Brown M., Heller J., (2003). Multivariate Non-Linear Feature Selection with Kernel Multiplicative Updates and Gram-Schmidt Relief. In *Proc. BISC FLINT-CIBI 2003 Workshop, Berkeley*.
5. Brown, M. P. (2000). Small subunit ribosomal rna modeling using stochastic context-free grammars. In *ISMB 2000*, pages 57–66.
6. Brown, M. P. S., Grundy, W. N., Lin, D., Cristianini, N., Sugnet, C., Furey, T. S., M.Ares, J., and Haussler, D. (2000). Knowledge-based analysis of microarray gene expression data using support vector machines. *PNAS*, 97(1):262–267.
7. Sjölander, K., Karplus, K., Brown, M., Hughey, R., Krogh, A., Mian, I., and Haussler, D. (1996). Dirichlet mixtures. *CABIOS*, 12(4):327–345.
8. Brown, M. P. S. and Wilson, C. (1995). Rna pseudoknot modeling using intersections of stochastic context free grammars with applications to database search. In Hunter, L. and Klein, T., editors, *Pacific Symposium on Biocomputing*, pages 109–125.
9. Sakakibara, Y., Brown, M., Hughey, R., Mian, I. S., Sjölander, K., Underwood, R. C., and Haussler, D. (1994a). Stochastic context-free grammars for tRNA modeling. *NAR*, 22:5112–5120.
10. Krogh, A., Brown, M., Mian, I. S., Sjölander, K., and Haussler, D. (1994). Hidden Markov models in computational biology: Applications to protein modeling. *JMB*, 235:1501–1531.
11. Sakakibara, Y., Brown, M., Hughey, R., Mian, I. S., Sjölander, K., Underwood, R. C., and Haussler, D. (1994b). Recent methods for RNA modeling using stochastic context-free grammars. In *Proc. Asilomar Conf. on Combinatorial Pattern Matching*, New York, NY. Springer-Verlag.
12. Brown, M. P., Hughey, R., Krogh, A., Mian, I. S., Sjölander, K., and Haussler, D. (1993). Using Dirichlet mixture priors to derive hidden Markov models for protein families. In Hunter, L., Searls, D., and Shavlik, J., editors, *First International Conference on Intelligent Systems for Molecular Biology*, pages 47–55, Menlo Park, CA. AAAI/MIT Press.

Presentations

Intelligent Systems for Molecular Biology 2000 :

Small subunit ribosomal rna modeling using stochastic context-free grammars.

Lawrence Berkeley National Laboratories 1999 :

SSU rRNA Modeling Using SCFGs and DNA Microarray Expression Data Analysis Using SVMs.

First Pacific Symposium on Biocomputing 1996 :

RNA pseudoknot modeling using intersections of stochastic context free grammars with applications to database search.

Software

RNACAD :

A system for modeling RNA using constrained SCFGs. Available through GPL at <http://www.cse.ucsc.edu/~mpbrown>.

The RNACAD system is a framework for modeling RNA secondary structure using SCFGs. It is able to estimate models from data, score new sequences, produce multiple alignments, and produce Fisher score vectors. It consists of a number of tools ranging from a top-down SCFG parser (C++), support scripts (Perl), and a visualization program (Java).

The RNACAD system is used at one of the premiere ribosomal RNA databases, the Ribosomal Database Project-II <http://rdp.cme.msu.edu/index.jsp>, that contains 190,785 aligned and annotated 16S rRNA sequences. All sequences are aligned against a general Bacterial rRNA alignment model using a modified version of the program RNACAD, a Stochastic Context Free Grammar (SCFG) based rRNA aligner that directly incorporates rRNA secondary structure information into its internal model.