

## Supplementary Material

# A Study and Analysis of Disease Identification using Genomic Sequence Processing Models: An Empirical Review

Sony K. Ahuja<sup>1</sup>, Deepti D. Shrimankar<sup>1,\*</sup> and Aditi R. Durge<sup>1</sup>

<sup>1</sup>Visvesvaraya National Institute of Technology, Computer Science and Engineering, India

### DATASETS AND ITS CHARACTERISTICS

SR NO	DATASET	DISEASE TYPE	CHARACTERISTICS
1.	<p>The Cancer Genome Atlas Program (TCGA)</p> <p>The Cancer Genome Atlas (TCGA), a landmark cancer genomics program, molecularly characterized over 20,000 primary cancer and matched normal samples spanning 33 cancer types.</p>	<p>Ovarian Cancer</p> <p>Colorectal Cancer</p> <p>Leukaemia</p> <p>Breast Cancer</p>	<p>Gene expression and gene mutation data are collected from TCGA. The gene expression data sets (Agilent 244K Custom Gene Expression G450) and gene mutation data sets are downloaded using the TCGA2STAT package. (Version: March 2, 2016).</p> <p>Gene expression profiles of colorectal cancer patients is collected from the TCGA database and statistical analysis is performed based on the data. Three clinical variables (stage, gender and age) are compared with gene expression profiles, and it is found that stage was the most influential variables on gene expression.</p> <p>Algorithm is applied to DNA methylation and gene expression data from The Cancer Genome Atlas (TCGA) Acute Myeloid Leukemia cohort to predict survival.</p> <p>Gene expression data sets of breast cancer is collected from The Cancer Genome Atlas which quantify the activity of genes from different cancer stages. There are four main progression stages, i to iv, of breast cancer included in the TCGA data sets. The gene expression levels of 17280 genes (Affymetrix HT Human Genome U133 Array Plate Set) for 496 breast cancers were used.</p>
2.	<p>Long noncoding RNAs (lncRNAs) database.</p> <p>The LncRNADisease database is not only a resource that curated the experimentally supported lncRNA-disease association data but also a platform that integrated tool(s) for predicting novel lncRNA-disease associations.</p>	Lung cancer	<p>The latest version of lncRNA-disease associations was downloaded from the LncRNADisease database. After removing the duplicate data, a gold-standard dataset was obtained, which included 855 experimentally validated LncRNADisease associations. (including 472 lncRNAs and 130 diseases)</p>



		Leukaemia	Both the training and test datasets include pairs of contrasted biomolecular networks along with the gene expression data for each gene. The expression data of each gene were collected from databases, such as the Gene Expression Omnibus (GEO) ( <a href="https://www.ncbi.nlm.nih.gov/geo/">https:// www.ncbi.nlm.nih.gov/geo/</a> ).
		Asthma	Two asthma gene expression datasets are used in the experiments, which are downloaded from the gene expression omnibus (GEO) repository. The first dataset (ID: GSE43696) contains 88 disease samples and 20 control samples, and each sample includes the expression data of 9194 genes. The second dataset (ID: GSE31773) contains 24 disease samples and 16 control samples, and each sample includes the expression data of 8789 genes.
		Multiple Disease types	Five SNP microarray datasets were used which are publicly available from NCBI GEO. GSE67047 series includes SNP data from patients who have sporadic Medullary Thyroid Cancer (sMTC) and juvenile Papillary Thyroid Cancer (PTC). There are 225 individuals in the data set, 96 patients and 129 parents, and approximately 1,000,000 features. GSE9222 series includes SNP data associated with Autism. The dataset consists of 567 individuals, 335 patients and 232 parents, and more than 250,000 features. GSE34678 series includes SNP data from patients who have Colorectal Cancer. This study utilizes 124 samples, 62 as a case and 62 as a control and approximately 250,000 features. GSE13117 series includes SNP data associated with Mental Retardation taken from affected patients and their healthy parents. The dataset consists of 360 individuals, 120 children and 240 parents, and nearly 250,000 features. GSE16619 series includes SNP data related to Breast Cancer with more than 500,000 SNPs. This study used 111 individuals, 69 as cases and 42 as controls.

5.	<p>NCBI Virus Database</p> <p>NCBI Virus database is an integrative, value-added resource supporting retrieval, display and analysis of a curated collection of virus sequences and large sequence datasets.</p>	Viruses	<p>Virus genomic sequences were retrieved from publicly available databases. In particular, five datasets, each representing one virus type, were considered. These datasets represent the viruses of Dengue, Hepatitis B, Hepatitis C, HIV-1, and Influenza A. The Dengue virus sequences were downloaded from the National Center for Biotechnology Information (NCBI) virus database. The same database was the source of Influenza genomes. The genomic sequences of Hepatitis B were downloaded from The Hepatitis B Virus Database (HBVdb) while Hepatitis C genomic sequences came from the database of Los Alamos National Laboratory (LANL). The latter database was also the source of HIV-1 genomes.</p>
6.	<p>The database of Genotypes and Phenotypes (dbGap)</p> <p>The database of Genotypes and Phenotypes (dbGaP) was developed to archive and distribute the data and results from studies that have investigated the interaction of genotype and phenotype in Humans.</p>	Pre-term birth classification	<p>The data, for this study, was obtained through authorised access to dbGap (Study Accession: phs000332.v3.p2). The dataset includes 722 cases and 1,057 controls. Cases were drawn from deliveries at the Boston Medical Center (BMC) that occurred before 37 weeks of gestation irrespective of birth weight. Controls include mothers who delivered term babies after 37 weeks of gestation also from the BMC cohort.</p>
7.	<p>National Taiwan University Hospital and the National Cheng Kung University Hospital.</p>	Bipolar disorder BD-I and BD-II	<p>The participants were 316 Han Chinese patients (195 with BD-I and 121 with BD-II) recruited from outpatient and inpatient settings</p>
8.	<p>The Hypertension (HT) dataset of the Wellcome Trust Case-Control Consortium (WTCCC)</p> <p>The Wellcome Trust Case Control Consortium (abbreviated WTCCC) is a collaboration between fifty research groups in the United Kingdom in the field of human genetics. Established in 2005, the WTCCC aims to conduct genome-wide association studies (GWASs) to shed light on the genetic architecture of common human diseases.</p>	Hypertension	<p>2,001 samples from individuals with hypertension disease were retrieved from the Hypertension (HT) dataset of the Wellcome Trust Case-Control Consortium (WTCCC), while 1,500 samples of control (healthy) individuals were retrieved from the UK National Blood Service Control Group (NBS)</p>

9.	GSE5281 brain dataset	Alzheimer's disease	The dataset contains 161 samples out of which 87 are diagnosed with Alzheimer's disease and 74 samples are from the healthy control group.
10.	Alzheimer's Disease Neuroimaging Initiative (ADNI) database.  ADNI is a multisite study that aims to improve clinical trials for the prevention and treatment of Alzheimer's disease.	Alzheimer's disease	Data was gathered from the ADNI database. ADNI Phase 1 data were collected from subjects (214 Non-AD and 177 AD cases). ADNI-1 samples were genotyped using the Human 610-Quad BeadChip, resulting in 620,901 SNPs.
11.	The University of California Irvine (UCI) repository.  The UCI Repository is a collection of real-life data, domain theories, and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms.	Heart Disease	The UCI repository was used to retrieve the heart disease database. The original database contains 76 attributes, but based on extensive experiments it was found that the most effective attributes were 14. The database contains the most dominant 14 attributes, which why they were chosen for training the model.