

1 *WDR88*, *CCDC11*, and *ARPP21* genes indulge profoundly in the
2 desmoplastic retort to prostate and breast cancer metastasis.

3 **Rajni Jaiswal^{1,*}, Shaurya Jauhari² and S.A.M. Rizvi²**

4 ¹**Department of Biotechnology, School of Engineering and Technology,**
5 **Sharda University, Greater Noida 201306.**

6 ²**Department of Computer Science, Jamia Millia Islamia, New Delhi 110025.**
7

8 **Abstract.** Microarray technology has unlocked doors to a multitude of open analysis prob-
9 lems that if conceived with efficacy may uncover varied genotypic and phenotypic traits. Al-
10 gorithms belonging to different cultures in computer science have been applied to gene ex-
11 pression data to derive correlation and stratification parameters. While most outcomes are
12 subject to clinical validation, majority of which get declined, the search for the precisely tar-
13 geted therapeutic agents is still on. This paper is an effort in the similar direction and strives
14 to delineate genes with significant stromal signatures. We suggest a corroborative indulgence
15 of a human laterality disorder gene, *CCDC11* in the metastasis, in addition to the role of
16 *WDR88* and *ARPP21* genes has been further materialized in the analysis. Another standout
17 aspect of the study has been the associated implications of the genes in rare disorders of male
18 breast and female prostate cancers. There is also a threshold proposal that stratifies “safe” ex-
19 pression space for genes. Complimentarily, the manuscript serves as an expedient protocol
20 for anyone seeking microarray data analysis, particularly in R.

21
22 **Keywords:** Breast Cancer, Desmoplasia, Gene Expression Data, Human Laterality Disor-
23 der, Microarray Analysis, Prostate Cancer, Reactive Stroma.

¹ Corresponding Author. Email: irajnijaiswal@gmail.com, Phone: +91 8744902881.
Article Type: Research Article, Manuscript Length: 7651, Tables: 4, Figures: 13

2

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

24 **Abbreviations:**

25 ANN Artificial Neural Networks

26 AR Androgen Receptor

27 DEG Differentially Expressed Genes

28 ECM Extracellular Matrix

29 ER Estrogen Receptor

30 FDA Fisher Discriminant Analysis

31 GRN Gene Regulatory Network

32 GSEA Gene Set Enrichment Analysis

33 GWAS Genome Wide Association Studies

34 IHC Immunohistochemistry

35 LCM Laser Capture Microdissection

36 NGS Next Generation Sequencing

37 PCA Principal Component Analysis

38 PCR Polymerase Chain Reaction

39 PSO Particle Swarm Optimization

40 SVM Support Vector Machines

41 HER2 Human Epidermal growth factor Receptor 2

42 PR Progesterone Receptor

43

44

45 **1. Introduction**

46 There appears nothing proverbially eerie about the technology at the get go. Microarrays
47 usage thrive with (Schena et al. 1995) and were originally applied to harbour global gene ex-
48 pression (DeRisi et al. 1997; DeRisi et al. 1996) in association with yeast studies. With the
49 proliferation of data pertaining to medication and that too in the digital proforma, it is crucial
50 to constantly challenge and update the current configuration of systems that are being used to
51 analyse it [genomic data] for compliance to the medical care. NGS is one such advancement
52 that was gullible to the geneticists. Unlike the microarray data that catalogues gene expres-
53 sion values under a predefined probe, the RNA-seq data from NGS documents expression
54 range in totality (Uziela & Honkela 2013). RNA sequencing technology pictures a compre-
55 hensive view of the transcriptome with the data being reproducible for novel discoveries
56 yielded by disparate analyses. RNA sequencing is also helpful in detection of structural varia-
57 tions as gene fusions, alternative splicing events, etc. But microarrays still continue to pro-
58 vide a relatively affordable *first-foot* to genomics, bearing robustness and short turn-around
59 time. With significant disparities owing to the definite and specific backgrounds of the indi-
60 viduals, the genomic data available via microarray format has shown likewise results when
61 particular maladies come into question as cancer, diabetes, amongst others. The big question
62 however is that could the genes be standardized via ontology driven mechanism so that spe-
63 cific drug targets be known and hunted for. Scientists are always looking for particular bi-
64 omarkers that can be universally acclaimed and acknowledged. In the current paper, we at-
65 tempt to underline key players that are actively responsible for representing the metastatic
66 behaviour and proliferation of oncogenic state in a body induced with breast-type and pros-

4

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.
67 tate-type cancers, in cognizance to stromal reaction. The results are based on a comparative
68 meta-analysis.

69 Reactive stroma is a response to the aberration into the tissues due to tumor invasion
70 (Planche et al. 2011). Synonymous to desmoplasia, it has also been recognized that the stro-
71 mal response is exclusive to tumor type. It can be perceived that desmoplastic response is in
72 tandem to carcinogenesis and subsequent metastasis. Thus, it is unstated that desmoplastic re-
73 action is also a prospective antecedent of premalignant stage, as the growth of connective, fi-
74 brous tissues around the tumor cells commences. Genetic irregularity in the cells compart-
75 mentalized in epithelium represents carcinoma in situ and the lesions initiate cell fibroblasts
76 as a tackling measure. Functionalities of stromal initiation include homeostasis and tissue
77 structure restoration. Chronicled is also that cell division govern mechanism is hampered be-
78 cause of the tumor induction and eventual progression. The amount of reactive stroma is pro-
79 portional to the disease state (Martin & Rowley 2013). Once the tumor foray infiltrates
80 through the ECM into adjoining host tissues, they become potent to further metastasize. Vas-
81 cular structures, blood and lymph vessels, ECM, and fibroblasts constitute the stroma (Casey
82 et al. 2009). Diverse studies by (Tuxhorn et al. 2002), (Ayala et al. 2003), (Roepman et al.
83 2006), and (Finak et al. 2006) implemented LCM to scrutinize gene expression profiles of
84 tumor stroma (breast) versus normal epithelium and clinched that the alterations in the stro-
85 mal microenvironment is comparative to the tumor progression.

86 In the following work, we attempt to ascertain genes that are prominent to tumor progres-
87 sion and subsequent stromal response. This may aid identification of key pathways (genes in-
88 stituted) that are liable for the cancer metastasis. As the dataset may reveal, we attempt to an-
89 alyze breast and prostate oncogenes.

90 This paper is organized as follows:

91 First, the developments in the breast cancer and prostate research, over the years, are cata-
92 logued. Various data analysis methodologies that have inferred some very seminal results
93 have been underlined. We then present our viewpoint and improvements in the domain and
94 propose a novel algorithm to analyse cancer stroma data. As it would necessitate, the sifted
95 targets are subjected to validation; but due to accessibility constraints, could only be done via
96 available erstwhile published research work. Their [genes] analysis can further substantiate
97 our studies for preventing the spread of cancer to the other tissues through pathway blockage
98 and rendering them benign through a drug treatment.

99 The statistical analysis and visualizations are covered with R language (version 3.2.3) (in-
100 terface used is RStudio version 0.99.491) on a desktop computer with 8 Gb RAM and an Intel
101 i5 CPU with 3.50 Ghz clock speed. For further distillations, MeV version 4.9.0 (Anon n.d.),
102 and Cytoscape are employed.

103

104 **1.1 An Alarming Statistic**

105

106 A not so long ago article (Kamath et al. 2013), reports that India is overwhelmed with 2.5
107 million cancer patients in aggregate and close to a million such augmented annually. To put
108 things into perspective, there were 1.7 million and 11.4 million cancer incidences in the
109 South East Asia region and world over, respectively in 2004. According to Globocan data
110 (International Agency for Research on Cancer), India tops the chart with 1.85 million years of

6

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

111 healthy life lost due to breast cancer alone. The aftermaths of this malady are equally likely
112 for the rest of the world too.

113

114 * *Healthy life lost is defined by years lost owing to premature death and deterioration of*
115 *health standards on account of a disease induction into the body.*

116

117 An elucidation from an erstwhile research confirms that after cervical cancer, breast cancer
118 is highly promulgated amongst Indian women. Also shown is that Indian women are likely to
119 inhibit breast cancer, a decade earlier than their Western counterparts. The paucity of early
120 detection and incompetent control mechanism can largely explain the succumbing rate. Ex-
121 orbitance in breast cancer cases throughout developing nations is proportionate to varying
122 lifestyle being is unregulated and sporadic, expectancy and delivery of fewer children, and
123 hormonal intervention exemplified by post-menopausal hormonal therapy. The symptoms are
124 profound at a later stage of the malignancy and hence pose greater challenge to review the
125 disease at the initiation. The authors of the study (Kamath et al. 2013) stressed upon the need
126 to exorcise this “ticking time bomb” and called for apt administrative measures for the same.

127 Prostate cancer, mostly occurring in elderly men, has similar danger trail and accounts for
128 second largest cancer causing deaths in U.S. males after lung cancer (Siegel et al. 2016;
129 Gaylis et al. 2016).

130

131

132 **1.2 Provenance**

133 When it comes to being most defiant and stubborn, and not to mention “incurable”, cancer
134 is christened far and wide for being the malady that poses serious threat to the manhood.
135 Many of the responsible genes involved in the pathways oriented to oncological disorders
136 have complex and overlapped functioning. Not to mention, some genes remain dormant at an
137 instance and are activated by a particular range of expression level of other corresponding
138 gene[s]. They also tend to become chemotherapy resistant through a self-regulatory mecha-
139 nism. These attributes account for a thorough and complacent inspection of the various pa-
140 rameters involved in gene functioning, mapped and homed-in.

141 In the exploration of gene expression data, the magnitude of tissue samples is lower with
142 respect to number of genes that may inevitably lead to overfitting of data and inappropriate
143 results (Shen et al. 2007). Gene selection is critical to elucidate tissue classification as well as
144 to model complex genetic and molecular underpinnings, which explain the relation between
145 genes and varied biological phenomenon. The stability of the analysis model can be accom-
146 plished through it.

147 BRCA1 and BRCA2 are vehemently recognized for hereditary breast and ovarian cancer
148 proneness. These are human genes that produce tumor suppressor proteins that implicitly ini-
149 tiate DNA repair mechanism. If a mutation is detected in any of the aforementioned genes,
150 the susceptibility to espousing tumor inception is high (Anon n.d.). BRCA mutation may lead
151 to the following probabilities:

152

- 153 • *40%-80% for breast cancer*

8

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

- 154 • 11%-40% for ovarian cancer
- 155 • 1%-10% for male breast cancer
- 156 • *Up to 39% for prostate cancer*
- 157 • 1%-7% for pancreatic cancer

158

159 Likewise, if any other relative cancer genes could be deciphered by comprehensively ana-
160 lyzing the gene expression data and establishing their helm in metabolism via clinical valida-
161 tion, we can get closer to disease treatment and increased understanding towards biology.
162 (Bosdet et al. 2013) take BRCA mutation testing to a whole new level by incorporating the
163 Second Generation Sequencing and Third Generation Sequencing procedures, collectively
164 known as NGS, to deal with increasing number of tests that the people are willing to take to
165 judge their cancer proneness. This era of NGS renders reduced cost, greater efficiency and
166 high throughput. The assay defined uses automated small amplicon PCR followed by sample
167 pooling and sequencing with a second-generation instrument.

168

169 **1.3 Androgen Receptor: An observable *cause commune***

170

171 Classically abnormality in males associated with prostate cancer, androgen receptor re-
172 sponse has been apropos (Yu et al. 2000). AR gene isn't solely responsible to harbor design
173 and characteristic instructions for sex drive and hair growth, but also facilitate sexual physi-
174 ologies. Positioned on the long (q) wing of the X chromosome at the 12th position, the AR
175 gene encompasses cohorts of CAG repeat regions (*triplets* or *trinucleotide repeats*). The

176 strength of quantifiable occurrences of these DNA segments account for the proneness of the
177 prostate cancer and breast cancer; while some studies hold more repeats liable, others blame
178 lesser ones (Yu et al. 2000). Research also depicts that mutations in the AR gene are account-
179 able for prostate cancer instantiation (Nelson 2002) (Giovannucci et al. 1997), albeit somatic
180 in nature. In women, longer CAG repeats and polymorphisms may increase the risk of endo-
181 metrial and breast cancers (Mehdipour, Pirouzpanah, Kheirollahi, & Atri, 2010).

182

183 **1.4 Gene Selection**

184

185 While holding candescence to the fact that intergenic regions relegated as “junk DNA”
186 have long been undermined, numerous follow up studies have unraveled that non-coding
187 RNAs, amongst other “dark” regions have a profound effect on regulation of gene expression
188 (Birney et al. 2007) (Carninci et al. 2005) (Cheng et al. 2005) (He et al. 2008). Since microar-
189 rays are designed to study gene measurements, the aforementioned parameters are left dilut-
190 ed. This aspect holds its vitality and is sure to influence the end result. Notwithstanding, it
191 has been known that Particle Swarm Optimization Technique (PSO) has been meticulously
192 significant in harnessing gene selection (Shen et al. 2007) (Yuan & Chu 2007) (Shen et al.
193 2008) (Chuang et al. 2008) (Lin et al. 2008). Other approaches include Artificial Neural Net-
194 works (ANN) and Fisher Discriminant Analysis (FDA), to name a few. An ensemble meth-
195 odology involving Particle Swarm Optimization (PSO) and Support Vector Machines (SVM)
196 has been observed to be particularly critical to feature selection and cornering genes of inter-
197 est (Yeung et al. 2009).

10

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

198 **1.5 Elucidation of Cancer Subtypes**

199

200 Breast cancer is a neoplasm that with distinct subtypes has differently representable histo-
201 pathological features and response to systemic therapies (Dai et al. 2016). Patient age, tumor
202 size, and axillary lymph node status have been deciding factors as well (Schnitt 2010). Im-
203 munohistochemistry (IHC) biomarkers have been classically deployed to ascertain subtyping.
204 They entail Estrogen Receptor (ER), Progesterone Receptor (PR), Androgen Receptor (AR),
205 and Human Epidermal growth factor Receptor 2 (HER2). Back in the 70's , there were two
206 subtypes that became known to us, viz. (luminal epithelial) ER+ and ER- (Perou et al. 2000)
207 (Sorlie et al. 2003) (Alexe et al. 2007). Triple negative breast cancer is characterized by a
208 cancer subtype devoid of ER, PR, and HER2 gene expressions. Compounds like tamoxifen
209 (for ER), and trastuzumab (for HER2), are tactless in dealing with triple-negative breast can-
210 cer. It is chemotherapeutically challenging as it warrants a grouping of disparately rated drugs
211 to target each of the receptor. Owing to its profile, triple negative breast cancer is revered as a
212 basal-type. Another recent study (Vici et al. 2015), illumines reasonableness of the triple posi-
213 tive breast cancer.

214 From prostate cancer viewpoint, gene fusions between TMPRSS2 and ETS hierarchies
215 have been stressfully documented (Tomlins et al. 2006), and also with ERG genealogy
216 (Penney et al. 2016). Expression levels of genes *MUC1* and *AZGP1* were also shown to cate-
217 gorically underline exclusive subtypes of prostate cancer from clinicopathological stance
218 (Lapointe et al. 2004).

219 (Herschkowitz et al. 2007) orchestrated a pioneering work that led to elucidation of a novel
220 sub-type pertaining to breast cancer disorder. This new subtype, referred to as Claudin-Low
221 was implicit of low expression genes. Also, traditionally, tumor types could be classified as
222 basal epithelial-like group (ERs), an *ERBB2*-overexpressing group, and a normal breast-like
223 group (Davidson & Liu. 2010). Another feature discovery from a study by (Sorlie et al. 2001)
224 had confirmed the possible subdivision of the ER+ tumor type into two clusters with distinc-
225 tive gene sets having particular corresponding clinical outcome.

226

227 **2. Results and Discussion**

228 The exegesis is premeditated so as to elucidate a quantifiable threshold that stratifies
229 gene expression space in conjunction to normal and cancer stromal response states. We delib-
230 erate to identify key transcriptional features that determine the high dimensional feature space
231 and visualize their inter-linkups via a regulatory network illustration. This is always compli-
232 mentary to ascertain our knowledge about genes and their pathway-occurrence motifs. . .
233 (Figure 1)

234

235

236

237

238

239

240

241

242

12

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

243 **2.1 Anonymous Genes/ Probes**

244

245 We identified 27236 entries, while scrutinizing the annotation fields of the dataset that
246 eluded ontological reference. There are also considerable amount of genes whose expression
247 values are catalogued under incongruent probes, resulting in their multiple incidences. This is
248 a purported case of genes' splice variants, as the dataset suggests. While we aim to identify
249 DEG and construct a respective GRN, there is also a prudence of elucidating functionally co-
250 herent genes that may unravel profiling of all or few anonymous genes. Thus, it appears du-
251 teous to abandon blank values to sustain quality of biological interpretation and germaneness.

252

253 **2.2 Normality and Data stabilization**

254

255 The data appears normalized data sans log transformation. Hence it is log-transformed and
256 metamorphosed to render mean=0, and standard deviation=1, i.e. it followed normal distribu-
257 tion. Since, the normality isn't skewed as a result of multiple comparisons problem (Dunn
258 1961), as we're not envisioning multitude of significance values, there is no proliferation of
259 Type I error occurrence anomalies. . . (Figure 2)

260

261

262

263

264

265

266 **2.3 Differentially Expressed Genes**

267

268 With respect to the assumed significance level (α) to be 0.01 and hence a stern confidence
269 interval of 99%, we aim to copiously optimize the gene(s) search by postulating as follows:

270

271 *(Null hypothesis)H₀: Genes are not differentially expressed (equal means)*

272 *(Alternate hypothesis)H₁: Genes are differentially expressed (unequal means)*

273

274 The listing of differentially expressed genes will implicitly catalogue up- and down-
275 regulated genes too. To prudently list them out, a within genes correlation does the job. The
276 negative numbers represent down-regulated genes and positives up-regulated ones.
277 (Danielsson et al. 2013) report that maximum of genes en route malignancy, are down regu-
278 lated. This is not for reference, but only to mark. There is also to note that since breast cancer
279 and prostate cancer find unique origins pertaining sex discrepancy, it's only logical to work
280 with bifurcated dataset. We contemplate breast cancer and prostate cancer entries with dis-
281 tinct exegesis and later combine and compare the results owing to significance to biological
282 interpretation.

283 The exploration renders 356 probes being differentially expressed in breast stroma and 221
284 in prostate stroma (with p-value < 0.01) amongst which *ADH1B*, *COL10A1* are most distinct-
285 ly expressed in breast strata, while *BMP5*, *SFRP4* are notable enough in prostate cluster. . .
286 (Figure 3) . . .

287 (Figure 4)

14

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

288 . . . (Figure 5)

289 . . . (Table 2)

290

291 We also acknowledge that packages like *siggenes* (Schwender, 2012), *samr* are available
292 that incorporate Significance Analysis of Microarrays (SAM) (Tusher, Tibshirani and Chu,
293 2001) working procedure, but in view of keeping the analysis more abstract and interactive,
294 there is a minimum use of readymade library functions.

295

296 **2.4 Gene Set Enrichment Analysis**

297

298 . . . (Figure 6)

299

300 Gene Set Enrichment Analysis (GSEA) is a scheme to map statistically relevant genes to
301 pre-known biological profiles, eg. phenotype, to discern their life relevance. The molecular
302 signatures are updated as the curation cascades. There exist a consortium of metadata librar-
303 ies for cataloguing genes and gene products' information. To standardize the practice of an-
304 notation in genomics, this bioinformatics initiative is absolutely imperative as we're riding
305 the snowball of discoveries in GWAS. (In R language, GWAS is facilitated by Fischer's ex-
306 act test.)

307 This method is applied to the resultant set of differentially expressed genes and only those
308 with a reference in MSigDB (Subramanian, Tamayo, et al., 2005) (with a valid Unigene_ID,
309 Entrez_ID, and GO_TERM) are selected for further analysis. To accomplish the same, MeV

310 and Cytoscape (in parts) are used. The GO listing wasn't available for 10 probes in prostate
311 data and 11 in breast data, which led to their discard. At this stage, our dataset has 345 and
312 211 rows in breast and prostate data, respectively.

313

314 **2.5 Functionally Coherent Genes**

315

316 An important aspect that escorts investigation of the differentially expressed genes is the
317 strength of associativity between their tumor and normal roles. This can be explored using
318 correlation technique of statistical descent. Commonly known Pearson's product-moment (or
319 simply, correlation) coefficient helps establish connect between two linearly distributed vari-
320 ables. In simplified terms, Spearman's coefficient is a non-parametric version of Pearson's
321 coefficient with ranked data (Hauke & Tomasz 2011). Since, our dataset selection is so, we
322 would prefer using Pearson's correlation measure as opposed to Spearman's or Kendall's
323 which is equally effective (or more) for the *qualitative* data. Kendall's τ is based on concord-
324 ance and discordance. The question is to establish similarity between two distinct genes,
325 technically two expression vectors (Saeed et al. 2003). An expression vector spans expression
326 values vide all featured experimental conditions. Albeit microarrays are not known to cater
327 isoform expression detection as they are not absolute exhibitors of gene expression and rather
328 give a relative value (log ratios of hybridization intensities).

329 The Pearson's correlation coefficient (r) for class labels X and Y is mathematically repre-
330 sented as follows:

331

16

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

332
$$r = \frac{\sum(X-\bar{X})(Y-\bar{Y})}{\sqrt{\sum(X-\bar{X})^2}\sqrt{\sum(Y-\bar{Y})^2}}$$

333

334 The resultant genes in breast cancer and prostate cancer were tested for correlation
335 amongst themselves. With each gene confronting every other gene, a $\frac{n(n-1)}{2}$ comparisons are
336 expected in a pairwise matrix format. Since distance measure is hinged around mean values, a
337 mix of positive and negative integers is likely. It is to note here that notwithstanding the am-
338 plitude of correlation, there is a significance of signs in the correlation matrix. A negative (-)
339 number indicates that a gene is inhibiting another gene, while a positive (+) marks that the
340 two are expressing collaterally. To safeguard our conviction to the fullest, the gene list was
341 filtered with a dual-parameterized statistic. We sifted the genes with low p-value significance
342 and high correlation measure. The tables catalogue genes from cancer-duos, with correlation
343 > 0.95 and p-value $< 10^{-7}$.

344

345

346 **2.6 Gene Regulatory Networks (GRN)**

347

348 Transcriptional activity can be precisely monitored with GRNs (Chai et al. 2014). A visu-
349 alization of putative pathways and the absolute values that are symbolic of the degree of
350 strength between two components can bring out some very useful linkage information. After
351 the elucidation of differentially expressed genes, the inkling is to draw a correlation measure
352 amongst them to infer a relational matrix with values $\{-1, 0, 1\}$ with interpretations anti-
353 correlated, no dependence, and correlated, respectively. This notion is certainly not delimited

354 to “naivety of adjacency”. The informal theme is the distance measure, but logically it may
355 falsify the overall outcome due to inherent biases with the chip construction. Therefore, the
356 notion of correlated transcripts is revered more viable. A gene regulatory network is a visual-
357 ization of a set/part(s) of genes that result in myriad (all) of cell processes, including metabo-
358 lism, cell signaling and transduction, cell growth control etc., which is vital to understand the
359 dynamics of molecular biology (Karlebach & Shamir 2008). But, the mechanistic inference
360 of the architecture is subject to experimental biology, a wet lab gig (Davidson & Levin 2005).
361 Nevertheless, the disposability of GRNs can’t be disparaged as they provide a blueprint of the
362 underlying system and tellingly optimize our erudition.

363 Correlation establishes the linear propensity in-between variables. For pursuit of the same,
364 we deliberate a Bayesian approach. In continuum to our expedition with R, the packages,
365 BNArray (Chen et al. 2006), NATbox (Chavan et al. 2009) deploy probabilistic slant to de-
366 cipher gene interactions, where NATbox shows competitive proficiency (Chavan et al. 2009).
367 In this treatise, however, we’ve considered Cytoscape as a pliant tool for visualization the
368 transcriptional network in corroboration with the *GeneMania* plugin. The output network ma-
369 trix of genes was exported to Cytoscape for visualization and analysis.

370 Post validation of the transcriptional networks, gene **CCDC11**, which has traditionally
371 been revered for human laterality disorder (Perles et al. 2012; Narasimhan et al. 2015), has
372 been elicited to show strong propensity in both (breast and prostate) cancer profiles. The
373 **Coiled-Coil Domain Containing 11** or CCDC11 is a protein coding gene which is closely
374 associated with epidermis in amphibians and skin fibroblasts from *Homo sapiens*
375 (Narasimhan et al. 2015). Re-annotated as **Cilia and Flagella Associated Protein 53**
376 (**CFAP53**), the mutation in CCDC11 exhibits perturbed left-right asymmetry (Silva et al.

18

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

377 2016). As showcased in the particular analysis, it has thorough connectedness at the order of
378 11 (aggregate) to other genes, which may signify functional co-regulation. From the under-
379 standing, it is proposed that this hub gene could be responsible to stage the process of stromal
380 response and coordinate in the transcriptional activities of the same.

381 **WDR88** gene on chromosome 19, reversed WD repeat domain 88, is a protein-coding
382 gene and a branded marker for the onset of prostate cancer (Chinnaiyan et al. 2013). In a top-
383 ical finding, the gene has also been shown to have links with schizophrenia (Richards et al.
384 2016). 167 organisms have orthologs with human gene WDR88 that is conserved in chim-
385 panzee, Rhesus monkey, dog, cow, mouse, rat, chicken, and frog.

386

387 . . . (Figure 7)

388

389 Another gene **ARPP21**, located in chromosome 9, has been exceptionally highlighted in
390 the breast and prostate cancer profiles. It has been duly captured to be frequently deregulated
391 as is miR-128 (Pellagatti et al. 2010; Li et al. 2013). According to NCBI RefSeq (June 2012),
392 this gene encodes a cAMP-regulated phosphoprotein. The encoded protein is enriched in the
393 caudate nucleus and cerebellar cortex. A similar protein in mouse may be involved in regulat-
394 ing the effects of dopamine in the basal ganglia. Alternate splicing results in multiple tran-
395 script variants. It is thence fathomed that these hub genes could be responsible to stage the
396 process of stromal response and coordinate in the transcriptional activities of the same.

397

398 . . . (Figure 8)

399

400 From gene ontology, ARPP21 is also attributed to the response to stimulus, triggering cel-
401 lular response to heat; at any temperature higher than the optimal stimulus of that organism.

402 A deliberation to the current study also entails the exceptional, yet formidable idea of
403 cross-linkages of breast and prostate cancers. Although the rudiments of breast and prostate
404 are oriented towards females and males, respectively, nonetheless, an exceptional yet indeli-
405 ble facet of female prostate and male breast profiles has been dimly studied. According to the
406 American Cancer Society, breast cancer is aggregate 100 times less common in males than
407 females; that is to calibrate the lifetime risk of a male getting breast cancer is 1 in 1000. Con-
408 trastingly, the skene/ periurethral gland carcinoma (female prostate cancer, in generic terms)
409 is also found to contribute less than 0.0003 percent towards all genital cancers in women
410 (Dodson et al. 1994). The numbers aren't intellectually stimulating, albeit we choose to delve
411 a little deeper.

412

413 **2.7 Female Prostate and Male Breast Carcinomas**

414

415 Female prostate, i.e. Skene gland, named after Alexander Johnston Chalmers Skene,
416 who was a British gynecologist from Scotland, is a homologue for the male prostate organ
417 and its adenocarcinoma is a scarce occurrence. Elevated Prostate Specific Antigen (PSA) and
418 PSA Phosphatase (PSAP) are potent markers for detecting prostate cancer in general (female
419 as well as male prostate specimens). Owing to the rarity, the female prostate cancer isn't
420 thoroughly researched too. With the limited physiological understanding, an older case study
421 presented a female subject with advanced form of the disease. It was treated with convention-

20

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.
422 al surgery (Ueda et al. 2012), to eventually weed out all the spread. Other techniques includ-
423 ing radiotherapy (Korytko et al. 2012) have been sublimely effective as well.

424 The female prostate is acknowledged as a functional part of the urinary and reproductive
425 systems in female humans (Zaviacic et al. 2000). It is located on the anterior wall of the vagi-
426 na, around the lower end of urethra, on each side. A chance that skene gland could be a sec-
427 ondary cancer site is also plausible. Estrogen (Estradiol, Estriol, and Estrone) and Progester-
428 one are the two key enzymes/ hormones that regulate the female breast development,
429 menstrual cycle, and sexual function. They are also luminaries in the prostate region in the
430 female gerbils. Estrogen is present in both male humans and female humans, and can be
431 measured for analyzing cancer of the reproductive system subunits, viz. ovaries, testicles, etc.
432 The cancer of the Skene gland is also more recognized in older females, showing tangible le-
433 sions (Custodio et al. 2010). Additionally, it has also been extensively deliberated that a fami-
434 ly history of breast cancer and prostate cancer engenders augmented jeopardy to a postmeno-
435 pausal woman gestating breast cancer (Robinson et al. 2015), (Beebe-Dimmer et al. 2015).
436 The abscesses in the gerbils are also shown to be driven by progesterone. A case history of
437 multiple pregnancies and ageing could be culpable for the female prostate disorder (Oliveira
438 et al. 2011).

439 Owing to the limited case studies of Skene gland cancer, the symptoms of the disorder
440 aren't well acknowledged and etiology is apparently impervious. As general indications,
441 bleeding in the urethra, that could also accompanied by sporadic pain, are primarily contin-
442 gent to symptomatic treatments. If the following conditions hold, a quick visit to the physi-
443 cian is often advisable.

- 444 • Arduous, frequent, and often difficult urination

- 445 • Bleeding from the urethra
- 446 • Painful sexual intercourse and pubic area
- 447 • Erratic menstrual cycle

448

449 The causes for Skene gland cancer are diversely plethoric. They can include infection as
450 prostatitis, some sexually transmitted infections (STIs) as gonorrhea; Polycystic Ovarian
451 Syndrome (PCOS) that renders imbalance and frequently abundance of reproductive hor-
452 mones, cysts, and adenofibroma.

453

454 . . . (Figure 9)

455

456 Another malady, although uncommon but not to be belittled as the rate of occurrence in-
457 creases every year, is the Male Breast Cancer (MBC). Mainly, females are more vulnerable to
458 breast cancer, having stocky breast tissue; however, males have pertinent breast tissue as
459 well. Scientifically, mutated copies of *BRCA1* and *BRCA2* genes incubated by male humans
460 are proverbial causes for MBC. Tamoxifen and anti-hormonal drugs are FDA-approved
461 chemotherapeutics to treat breast cancer in both male and females. Requisite surgery (mastec-
462 tomy/ lumpectomy) followed by radiation therapy is standardly warranted, although individ-
463 ual therapies could include more aggressive treatment options. The ideology of a male human
464 incubating breast cancer is largely pondered with ignorance and aversion; this conviction, in
465 most cases, delays the screening of the disease. Peculiar symptoms of MBC entail ruptured
466 (and often painful) nipples, puckering and dimpled masses of the breast, decolorized jaggy

22

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

467 surfaces, etc. MBC is usually detected as a hard lump underneath the nipple and areola. The

468 histopathological derivatives in MBC and female breast cancer are homogenous.

469

470

471 . . . (Figure 10)

472 . . . (Figure 11)

473

474

475 Gynecomastia is also a disorder in men of benign nature, where the breast tissue becomes

476 enlarged due to hormonal misbalance (oestrogen to testosterone ratio), especially during pu-

477 berty. Although a natural phenomenon, it is usually conceived with humiliation and anxiety;

478 however paltry cases have been reported to establish that gynecomastia and MBC are con-

479 comitant. In conjunction, pseudogynecomastia is a condition when adipose tissue (fat) causes

480 gynecomastia.

481 Therefore, it can be argued that the denominations of origins, histopathology, causes,

482 symptoms, and treatments are overlapped for male-breast and female-breast cancer; and like-

483 ly so for male-prostate and female-prostate cancer. The contributing genes and pathways

484 could be further explored for overlap in disease profile and therapy.

485

486

487

24

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

509 where w is weight vector, x is input vector, and b is bias. The decision boundary can also
510 be deemed as a linear combination of support vectors. As calculated, 8 and 9 support vectors
511 were rendered from prostate and breast data respectively. The bias vectors from `<svm model-`
512 `el>$rho` are 1.429841 and 4.139861, from prostate and breast data correspondingly. Further
513 information can be found, as code output, from the supplementary documents.

514

515 ... (Figure 12)

516 ... (Figure 13)

517 ... (Table 3)

518 ... (Table 4)

519

520 **2.9 Conclusion and Future Work**

521 This text has been premeditated to render the most interactive portrayal of working with
522 gene expression data analysis. As a part of the original work, the authors have carried out
523 survival analysis too. The treatise however concentrates on the improved biomarker(s) dis-
524 semination.

525 As an imminent applicability, the study can aid fostering of pertinent therapeutics to deride
526 proliferation of cancer metastasis from one tissue to another by monitoring the expression
527 threshold and keeping it checked.

528 The procedure highlights an illustration of the packages available in the R language and
529 Bioconductor that duly facilitate the exploratory analysis of the genomic data. While doing
530 so, certain cohorts of genes were found relevant and were statistically narrowed to seed fur-
531 ther analysis. This aids reducing the search space for biomarkers (broadly explains the doc-

532 trine of bioinformatics) and the pipeline of wet laboratory testing and validation, proceeds. If
533 the genes *CDCC11*, *WDR88* and *ARPP21* have any causal implications in the stromal re-
534 sponse to the cancer metastasis, can and will only be substantiated through valid laboratory
535 studies. Researches alike add to the annotations of the known gene functionality. An array of
536 such explorations is warranted and is indeed happening. This trend over a period of time is
537 believed to pave way for a precision medicine schedule, when drug compounds' applications
538 and the respective gene functions are almost perfectly matched.

539

540

541 **3. Material and Methods**

542 **3.1 Dataset Selection**

543

544 The gene expression dataset chosen from the study is derived out of a study based on stro-
545 mal cells and invasive breast and prostate cancer development (Planche et al. 2011)

546

547 ... (Table 1)

548

549 The authors have commenced performing log transformation oriented normalization and
550 moved further with a primary cue gathering via Principal Component Analysis (PCA). It is
551 also reported that a very few number of overlain genes befall from breast and tumor profiles.
552 Pearson correlation coefficients exhibit stout propensity of breast stromal genes with breast

26

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.
553 data and prostate stromal genes with prostate data (Figure 1) (Planche et al. 2011). To add to
554 further consolidation of the outcome, survival analysis was carried out using Univariate Cox
555 approach that highlighted genes whose expression levels were crucially associated with the
556 patient survival. The downloaded dataset has no observable missing values in the cells to im-
557 pute; rather the blank entries are subsidiary to gene and probe ids.

558 Technically deduced from the background meta-analysis of the subject, we may decipher
559 that cancer will need a host medium (tissue) to proliferate to the other cells/ tissues/ organs.
560 The metastasis front of cancer would seek for the favorable restructuring of the basal tissue
561 framework. From the anticancer therapeutic vantage, hence, it renders incumbent that the on-
562 cogenes and stromal response must be equally thrust.

563 Through this exemplar multifaceted exegesis, we objectivize to construe the following:

564

- 565 a) Contrivance of differentially expressed genes (DEG)
- 566 b) GRN reconstruction, and
- 567 c) Decoding functionally coherent genes (eliciting anonymous genes) in accord to iso-
568 form expression.
- 569 d) Designing a classifier (machine learning approach) that embraces a threshold value
570 of gene expression that triggers ambient oncological desmoplastic response.

571

572 From statistical standpoint, the data concerned is *paired*, i.e. two different conditions (can-
573 cerous and normal, here) hybridized on the same slide. A recce exhibits noticeable gene en-
574 tries that outlie the tightly stratified expression space, as can be derived from the Fig. 3. The
575 dataset dimension of 54675 features tacitly conveys the infestation of multiple gene entries

576 associated with diverse probes. However, a cursory reconnaissance shall also establish that
577 there is only one replicate to each experimental condition.

578 In recent years, the molecular data has become reverently large. R has evolved as the *de-*
579 *facto* tool for genomic data analysis attributable to its IDE, flexibility and workflow control.
580 Amongst others Python is a viable option too. Biopython is a dedicated version of the lan-
581 guage for biological data analytics. However, R has an edge over other languages in terms of
582 packages (functionalities) to cope with the multidimensional data. Being open-source and ful-
583 ly distributable adds to the prowess as well.

584

585 **Declaration of Interest**

586 The authors register no conflict of interest.

587

588 **Author Contribution**

589 **Conception and design:** Rajni Jaiswal, Shaurya Jauhari.

590 **Development of methodology:** Rajni Jaiswal

591 **Analysis and interpretation of data (e.g., statistical analysis, biostatistics, computa-**
592 **tional analysis):** Rajni Jaiswal, Shaurya Jauhari.

593 **Writing, review, and/or revision of the manuscript:** Shaurya Jauhari, S.A.M. Rizvi.

594 **Study supervision:** S.A.M. Rizvi.

595

596

597

References

598 Alexe G, Dalgin G S, Ganesan S, DeLisi C and Bhanot G 2007 Analysis of breast cancer progression
599 using principal component analysis and clustering; *J. Biosci.* **32**, pp 1027–1039.

600

601 Birney E, Stamatoyannopoulos JA, Dutta A, Guigo R, Gingeras et al (2007) Identification and analysis
602 of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**: pp 799–
603 816.

604

605 Carninci P, Kasukawa T, Katayama S, Gough J, et al. (2005) The transcriptional landscape of the
606 mammalian genome, *Science* **309**: pp 1559–1563.

607

608 Casey T, Bond J, Tighe S, Hunter T, Lintault L, Patel O, Eneman J, Crocker A, White J, Tessitore J,
609 Stanley M, Harlow S, Weaver D, Muss H, Plaut K. Molecular signatures suggest a major role for stro-
610 mal cells in development of invasive breast cancer. *Breast Cancer Res Treat.* 2009 Mar; **114**(1): pp 47-
611 62. [doi: 10.1007/s10549-008-9982-8](https://doi.org/10.1007/s10549-008-9982-8)

612

613 Chai L, Loh S, Low S, Mohamad M, Deris S, Zakaria Z (2014), A review on the computational ap-
614 proaches for gene regulatory network construction, *Computers in Biology and Medicine*, Elsevier, 48
615 (1) : 55-65.

616

617 Chavan SS, Bauer MA, Scutari M, Nagarajan R, NATbox: a network analysis toolbox in R, *BMC Bio-*
618 *informatics* 2009, 10 (Suppl 11): S14, [doi: 10.1186/1471-2105-10-S11-S14](https://doi.org/10.1186/1471-2105-10-S11-S14)

619

620 Cheng J, Kapranov P, Drenkow J, Dike S, Brubaker S, Patel S, Long J, Stern D, Tammana H, Helt G,
621 Sementchenko V, Piccolboni A, Bekiranov S, Bailey DK, Ganesh M, Ghosh S, Bell I, Gerhard DS,

622 Gingeras TR (2005) Transcriptional maps of 10 human chromosomes at 5-nucleotide resolution. *Sci-*
623 *ence* **308**: pp 1149–1154.

624

625 Danielsson F., Marie Skogs M., et al., Majority of differentially expressed genes are down-regulated
626 during malignant transformation in a four-stage model, PNAS 2013 110 (17) 6853-6858; published
627 ahead of print April 8, 2013, [doi:10.1073/pnas.1216436110](https://doi.org/10.1073/pnas.1216436110)

628

629 David Meyer, Evgenia Dimitriadou, Kurt Hornik, Andreas Weingessel and Friedrich Leisch (2015).
630 e1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071),
631 TU Wien. R package version 1.6-7. <https://CRAN.R-project.org/package=e1071>

632

633 Davidson P, Liu J, 2010 Analysis of Microarray Data to Confirm Novel Subtype of Breast Cancer;
634 *Thesis*, Worcester Polytechnic Institute.
635 [doi: 10.1016/j.compbimed.2014.02.011](https://doi.org/10.1016/j.compbimed.2014.02.011)

636

637 Dunn OJ, Multiple Comparisons Among Means, *Journal of the American Statistical Association*, 1961,
638 56(293).

639

640 Edward Giovannucci, Meir J. Stampfer, Krishna Krithivas, Myles Brown, Adam Brufsky, James Tal-
641 cott, Charles H. Hennekens, and Philip W. Kantoff (1997) The CAG repeat within the androgen recep-
642 tor gene and its relationship to prostate cancer PNAS **94** (7) 3320-3323

643

644 Eric Davidson, Michael Levin, Gene regulatory networks, *PNAS* 2005 102 (14) 4935; published ahead
645 of print April 4, 2005, [doi:10.1073/pnas.0502024102](https://doi.org/10.1073/pnas.0502024102).

646

30

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

647 Finak G, Sadekova S, Pepin F, Hallett M, Meterissian S, Halwani F, Khetani K, Souleimanova M,
648 Zabolotny B, Omeroglu A, Park M (2006) Gene expression signatures of morphologically normal
649 breast tissue identify basal-like tumors. *2006 Breast Cancer Res* 8:R58

650

651 Fodale V, Pierobon M, Liotta L, Petricoin E., Mechanism of cell adaptation: when and how do cancer
652 cells develop chemoresistance? *Cancer J.* 2011 Mar-Apr; 17(2):89-95. [doi:](https://doi.org/10.1097/PPO.0b013e318212dd3d)
653 [10.1097/PPO.0b013e318212dd3d](https://doi.org/10.1097/PPO.0b013e318212dd3d).

654

655 Gustavo Ayala, Jennifer A. Tuxhorn, Thomas M. Wheeler, Anna Frolov, Peter T. Scardino, Makoto
656 Ohori, Marcus Wheeler, Jeffrey Spitler, and David R. Rowley, Reactive Stroma as a Predictor of Bio-
657 chemical-Free Recurrence in Prostate Cancer, *Clin Cancer Res* October 15, 2003 9; 4792

658

659 Guy Karlebach & Ron Shamir, Modelling and analysis of gene regulatory networks, *Nature Reviews*
660 *Molecular Cell Biology* 9, 770-780 (October 2008) | [doi: 10.1038/nrm2503](https://doi.org/10.1038/nrm2503).

661

662 Hastie T, Tibshirani R, Friedman J, The Elements of Statistical Learning Data Mining, Inference, and
663 Prediction, 2nd ed., *Springer Series in Statistics*, ISBN 978-0-387-84857-0, [doi :10.1007/978-0-387-](https://doi.org/10.1007/978-0-387-84858-7)
664 [84858-7](https://doi.org/10.1007/978-0-387-84858-7).

665

666 Hauke J., Kossowski T., Comparison of values of Pearson's and Spearman's correlation coefficient on
667 the same sets of data. *Quaestiones Geographicae* 30(2), Bogucki Wydawnictwo Naukowe, Poznań
668 2011, pp. 87–93. [doi :10.2478/v10117-011-0021-1](https://doi.org/10.2478/v10117-011-0021-1)

669

670 He Y, Vogelstein B, Velculescu VE, Papadopoulos N, Kinzler KW (2008) The antisense transcrip-
671 tomes of human cells. *Science* 322: pp 1855–1857.

672

673 Ian E. Bosdet, T. Roderick Docking, Yaron S. Butterfield, Andrew J. Mungall, Thomas Zeng, Robin J.
674 Coope, Erika Yorida, Katie Chow, Miruna Bala, Sean S. Young, Martin Hirst, Inanc Birol, Richard A.
675 Moore, Steven J. Jones, Marco A. Marra, Rob Holt, Aly Karsan (2013) A Clinically Validated Diag-
676 nostic Second-Generation Sequencing Assay for Detection of Hereditary BRCA1 and BRCA2 Muta-
677 tions, *The Journal of Molecular Diagnostics*; doi: [10.1016/j.jmoldx.2013.07.004](https://doi.org/10.1016/j.jmoldx.2013.07.004)
678
679 James G, Witten D, Hastie T, Tibshirani R, An Introduction to Statistical Learning with Applications in
680 R, *Springer Texts in Statistics*, ISBN 978-1-4614-7137-0 doi: [10.1007/978-1-4614-7138-7](https://doi.org/10.1007/978-1-4614-7138-7).
681
682 Jason I Herschkowitz, Karl Simin, Victor JWeigman, Igor Mikaelian, Jerry Usary, Zhiyuan Hu, Karen
683 E Rasmussen, Laundette P Jones, Shahin Assefnia, Subhashini Chandrasekharan, Michael G Backlund,
684 Yuzhi Yin, Andrey I Khramtsov, Roy Bastein, John Quackenbush, Robert I Glazer, Powel H Brown,
685 Jeffrey E Green, Levy Kopelovich, Priscilla A Furth, Juan P Palazzo, Olufunmilayo I Olopade, Philip S
686 Bernard, Gary A Churchill, Terry Van Dyke, Charles M Perou, (2007) Identification of conserved gene
687 expression features between murine mammary carcinoma models and human breast tumors, *Genome*
688 *Biology*, **8**: R76, doi: [10.1186/gb-2007-8-5-r76](https://doi.org/10.1186/gb-2007-8-5-r76)
689
690 Jennifer A. Tuxhorn, Gustavo E. Ayala, Megan J. Smith, et al, Reactive Stroma in Human Prostate
691 Cancer: Induction of Myofibroblast Phenotype and Extracellular Matrix Remodeling, *Clin Cancer Res*
692 2002; **8**: pp 2912-2923.
693
694 Kalluri R, Zeisberg M, Fibroblasts in cancer, *Nature Reviews Cancer* **6**, 392-401 (May
695 2006) | doi:[10.1038/nrc1877](https://doi.org/10.1038/nrc1877)
696
697 Li-Yeh Chuang, Hsueh-Wei Chang, Chung-Jui Tu, Cheng-Hong Yang, (2008), Improved binary PSO
698 for feature selection using gene expression data, Elsevier, *Computational Biology and Chemistry* **32**,
699 pp 29–38.

32

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

700

701 Martin, Rowley, Role of Reactive Stroma in Prostate Cancer, Prostate Cancer Protein Reviews, 2013,
702 **16**, pp 43-63.

703

704

705

706 Mehdipour, P., Pirouzpanah, S., Kheirollahi, M. and Atri, M. (2011), Androgen Receptor Gene CAG
707 Repeat Polymorphism and Breast Cancer Risk in Iranian Women: A Case-Control Study. The Breast
708 Journal, **17**: 39–46. doi: [10.1111/j.1524-4741.2010.01031.x](https://doi.org/10.1111/j.1524-4741.2010.01031.x)

709

710 National Cancer Institute, Factsheet, “BRCA1 and BRCA2: Cancer Risk and Genetic Testing”,
711 <http://www.cancer.gov/cancertopics/factsheet/Risk/BRCA>

712

713 Nelson K. A. and Witte J. S. (2002), Androgen Receptor CAG Repeats and Prostate Cancer, Am. J. Ep-
714 idemiol. **155** (10): 883-890. doi: [10.1093/aje/155.10.883](https://doi.org/10.1093/aje/155.10.883)

715

716 Nils Wilking et. al., “Prevention and the economic burden of the breast cancer”, Whitepaper, a report
717 commissioned by GE Healthcare, <http://newsroom.gehealthcare.com/>.

718

719 Perles Z et al., A human laterality disorder associated with recessive CCDC11 mutation, J. Med. Genet.
720 2012 Jun; 49(6): pp 386-390.

721

722 Perou CM, Sorlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, Pollack JR, Ross DT, Johnsen H,
723 Akslen LA, et al.: Molecular portraits of human breast tumours. *Nature* 2000, **406**: pp 747-752.

724

725 Pisco, A.O. et al. Non-Darwinian dynamics in therapy-induced cancer drug resistance. *Nature commu-*
726 *nications* **4**, 2467 (2013).

727

728 Planche A, Bacac M, Provero P, Fusco C, Delorenzi M, et al. (2011) Identification of Prognostic Mo-
729 lecular Features in the Reactive Stroma of Human Breast and Prostate Cancer. PLoS ONE **6**(5):
730 e18640. [doi:10.1371/journal.pone.0018640](https://doi.org/10.1371/journal.pone.0018640)

731

732 Qi Shen, Wei-Min Shi, Wei Kong, (2008) Hybrid particle swarm optimization and tabu search ap-
733 proach for selecting genes for tumor classification using gene expression data, Elsevier, Computational
734 Biology and Chemistry **32**, pp 53–60.

735

736 Qi Shen, Wei-Min Shi, Wei Kong, Bao-Xian Ye (2007) A combination of modified particle swarm op-
737 timization algorithm and support vector machine for gene selection and tumor classification. Talanta,
738 Volume 71, Issue 4, 15 March 2007, pp 1679-1683, ISSN 0039-9140, [doi](https://doi.org/10.1016/j.talanta.2006.07.047)
739 [:10.1016/j.talanta.2006.07.047](https://doi.org/10.1016/j.talanta.2006.07.047)

740

741 Roepman P, de Koning E, van Leenen D, de Weger R, Kummer JA, Slootweg P, Holstege F, Dissec-
742 tion of a metastatic gene expression signature into distinct components. 2006 *Genome Biol* **7**:R117

743

744 Saeed AI, Sharov V, White J, Li J, Liang W, Bhagabati N, Braisted J, Klapa M, Currier T, Thiagarajan
745 M, Sturn A, Snuffin M, Rezantsev A, Popov D, Ryltsov A, Kostukovich E, Borisovsky I, Liu Z,
746 Vinsavich A, Trush V, Quackenbush J. TM4: a free, open-source system for microarray data manage-
747 ment and analysis. *Biotechniques*. 2003 Feb; **34**(2):374-8.

748

749 Schena M., Shalon D., Davis R.W., and Brown P.O.: Quantitative monitoring of gene expression pat-
750 terns with a complementary DNA microarray. 1995 *Science* **270** (5235), pp 467-70.

751

752 Schwender H (2012). *siggenes: Multiple testing using SAM and Efron's empirical Bayes approaches*. R
753 package version 1.44.0.

Desmoplastic Retort to Prostate and Breast Carcinomas' Metastasis.

754

755 Sheng-Fa Yuana,b, Fu-Lei Chu, (2007) Fault diagnostics based on particle swarm optimization and
756 support vector machines, Sheng-Fa Yuan, Fu-Lei Chu, Elsevier, Mechanical Systems and Signal Pro-
757 cessing **21**, pp 1787–1798.

758

759 Shih-Wei Lin, Kuo-Ching Ying , Shih-Chieh Chen , Zne-Jung Lee (2008) Particle swarm optimization
760 for parameter determination and feature selection of support vector machines, Elsevier, Expert Systems
761 with Applications **35**, pp 1817–1824.

762

763 Significance Analysis of Microarrays (SAM), <http://statweb.stanford.edu/~tibs/SAM/>.

764

765 Sorlie T, Perou CM, Tibshirani R, Aas T, Geisler S, Johnsen H, Hastie T, Eisen MB, van de Rijn M,
766 Jeffrey SS, *et al.*: Gene expression patterns of breast carcinomas distinguish tumor subclasses with
767 clinical implications. *Proc Natl Acad Sci USA* 2001, **98**: pp 10869-10874.

768

769 Sorlie T, Tibshirani R, Parker J, Hastie T, Marron JS, Nobel A, Deng S, Johnsen H, Pesich R, Geisler
770 S, *et al.*: Repeated observation of breast tumor subtypes in independent gene expression datasets. *Proc*
771 *Natl Acad Sci USA* 2003, **100**: pp 8418-8423.

772

773 Subramanian, Tamayo, et al., Gene Set Enrichment Analysis, ([2005, PNAS 102, 15545-15550](#)),
774 Mootha, Lindgren, et al. ([2003, Nat Genet 34, 267-273](#)).

775

776 Subramanian, Tamayo, et al., Molecular Signatures Database (MSigDB), ([2005, PNAS 102, 15545-](#)
777 [15550](#))

778

779 Uziela K, Honkela A (2013) Probe Region Expression Estimation for RNA-Seq Data for Improved Mi-
780 croarray Comparability. PLoS ONE 10(5): e0126545. [doi:10.1371/journal.pone.0126545](https://doi.org/10.1371/journal.pone.0126545)

781

782 Walker R.A., The complexities of breast cancer desmoplasia, *Breast Cancer Res* 2001, **3**:143-
783 145, [doi:10.1186/bcr287](https://doi.org/10.1186/bcr287).

784

785 Whatcott CJ, Posner RG, Von Hoff DD, et al. Desmoplasia and chemoresistance in pancreatic cancer.
786 In: Grippo PJ, Munshi HG, editors. Pancreatic Cancer and Tumor Microenvironment. Trivandrum (In-
787 dia): Transworld Research Network; 2012. Chapter 8. Available from:
788 <http://www.ncbi.nlm.nih.gov/books/NBK98939/>

789

790 Xiaohui Chen, Ming Chen, and Kaida Ning, *BNArray*: an R package for constructing gene regulatory
791 networks from microarray data by using Bayesian network, *Bioinformatics* (2006) 22 (23): 2952-2954
792 first published online September 27, 2006 [doi:10.1093/bioinformatics/btl491](https://doi.org/10.1093/bioinformatics/btl491).

793

794 Yu H, Bharaj B, Vassilikos EJ, Giai M, and Diamandis EP (2000), Shorter CAG repeat length in the
795 androgen receptor gene is associated with more aggressive forms of breast cancer. *Breast Cancer Res*
796 *Treat.* **59**(2):153-61.

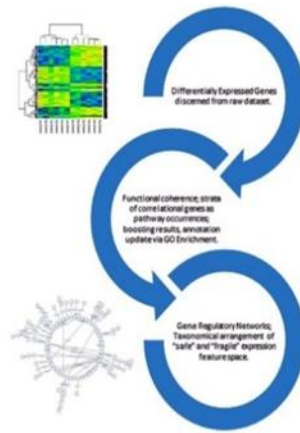


Figure 1 Illustration of workflow.

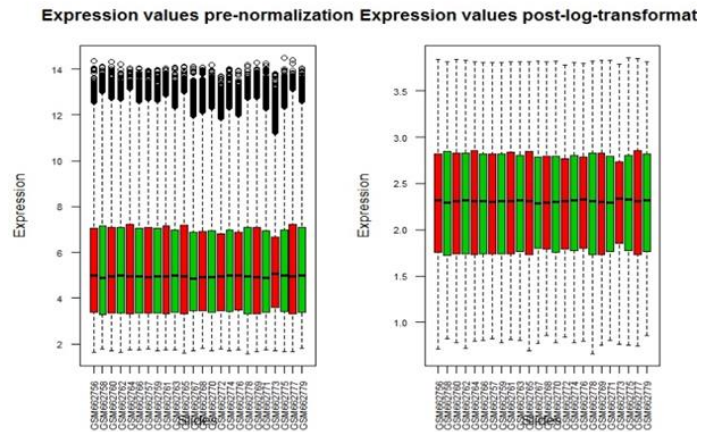


Figure 2 Box plots depicting sample expressions pre and post-normalization. Log₂ transformation is applied for the same and the data is rendered more balanced ahead of analysis. The cancer and non-cancer bars are represented by red and green color codes.

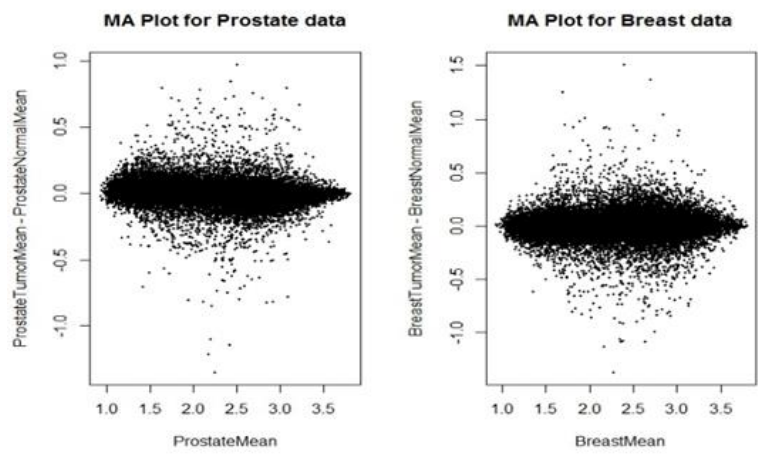


Figure 3 Respective MA plots of prostate and breast subsets. Clearly the floating specks demarcate the differentially expressed transcripts.

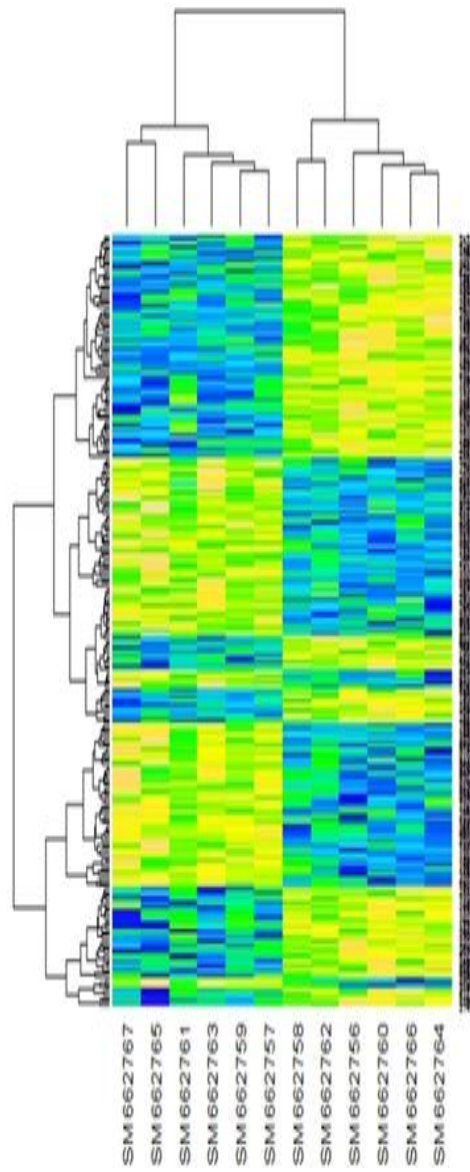


Figure 4 Heat map for differentially expressed breast genes. 356 probes with p -value < 0.01 were unraveled as being significant.

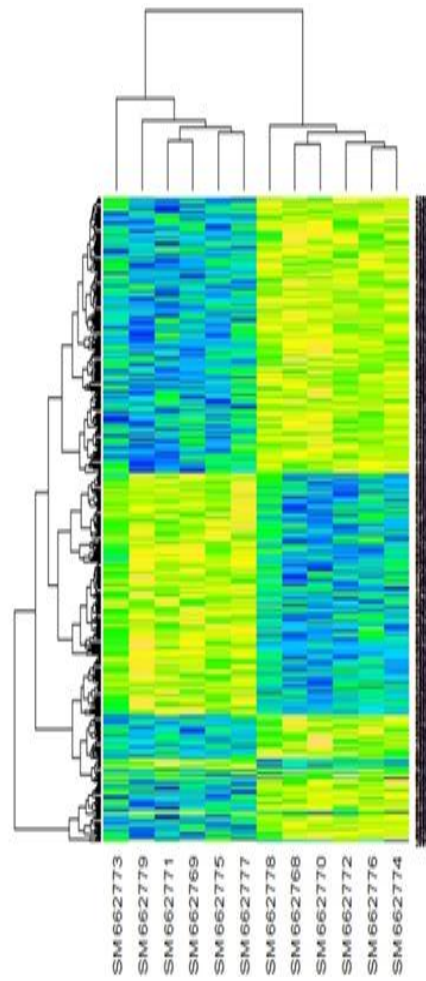


Figure 5 Heat map for differentially expressed prostate genes. Here, 221 probes with p-value < 0.01 were deemed crucial.

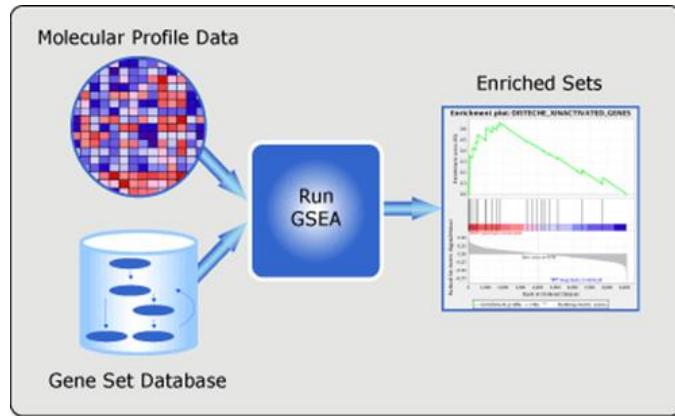


Figure 6 Illustration of GSEA framework (Subramanian, Tamayo, et al., 2005).

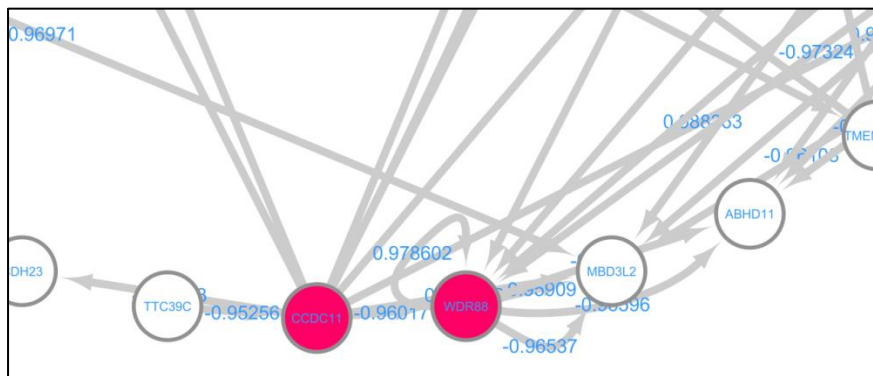


Figure 7 The study fortifies the eminent role of CCDC11 and WDR88 genes that are fundamental test genes for cancer diagnosis. The figure portrays a window from the prostate cancer GRN. As elicited, CCDC11 is orchestrating other genes, while WDR88 is a coveted

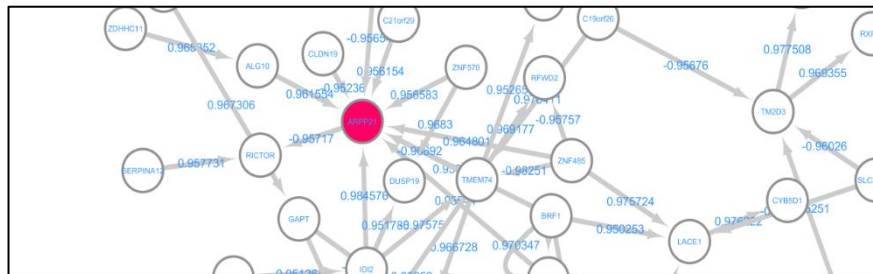


Figure 8 An excerpt from the Breast cancer GRN analysis shows profound coverage of ARPP21 gene with high propensity.

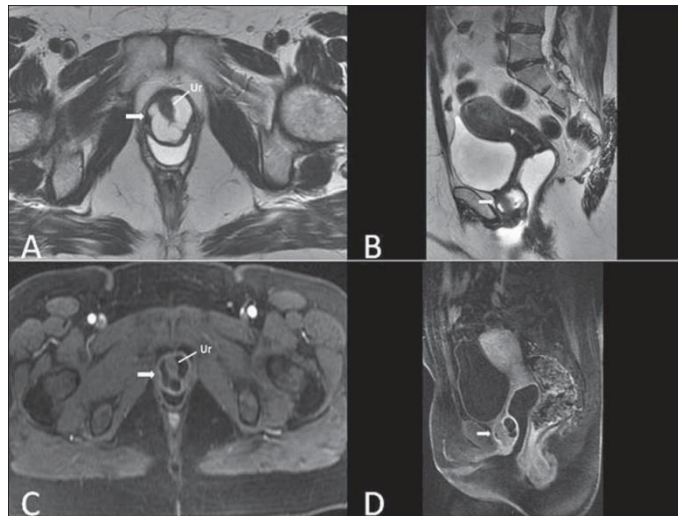


Figure 9 MRI scans of the cysts of the Skene glands. Multiplanar MRI T2-weighted (A,B) and contrast- enhanced T1-weighted (C,D) sequences identifying distal periurethral cysts (Ur) (arrows) located between the urethra and the vagina

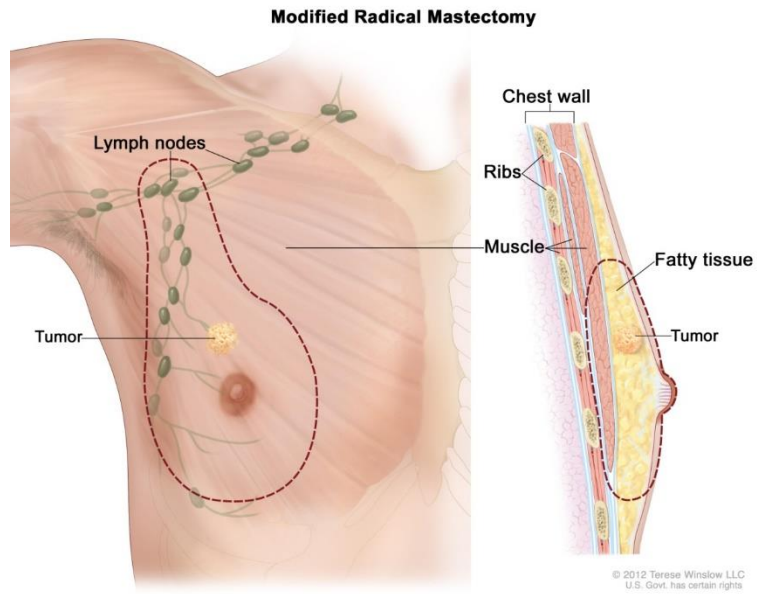


Figure 10 An illustration detailing radical mastectomy. Credit: <http://www.cancer.gov>

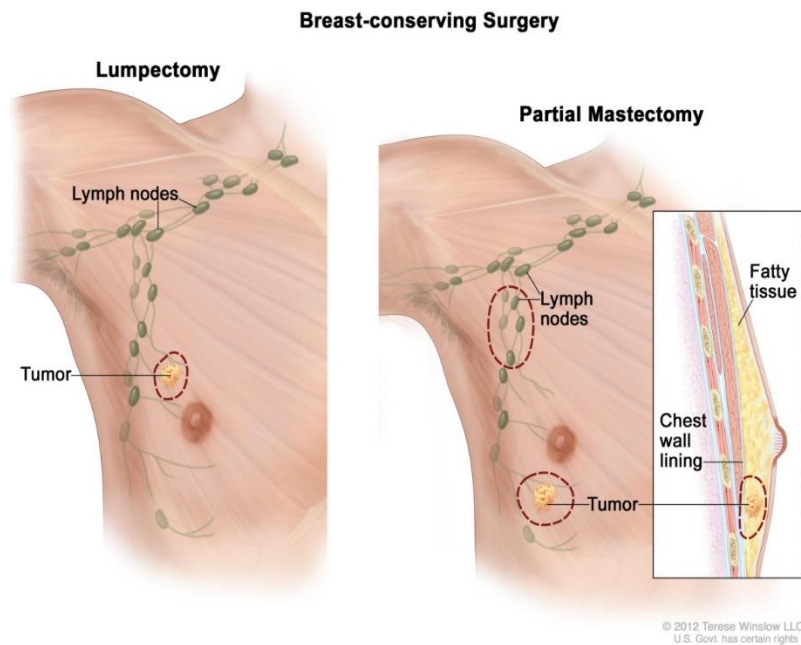


Figure 11 An illustration detailing breast-conserving surgery. Credit: <http://www.cancer.gov>

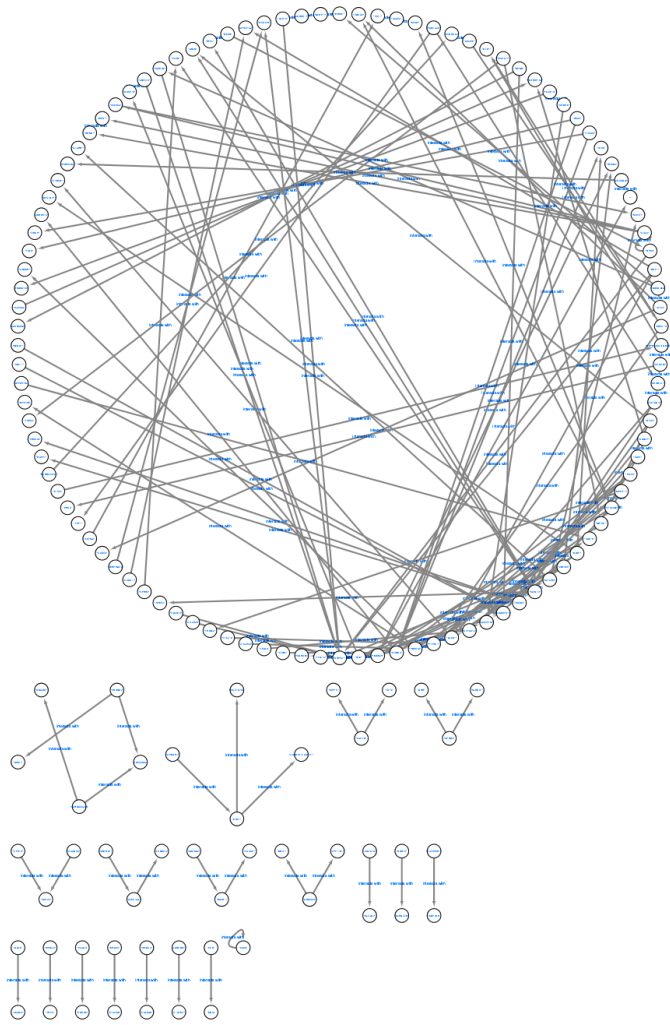


Figure 12 Breast Genes Regulatory Network

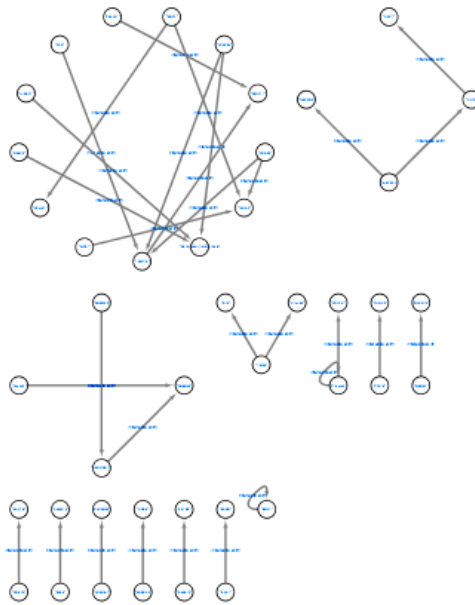
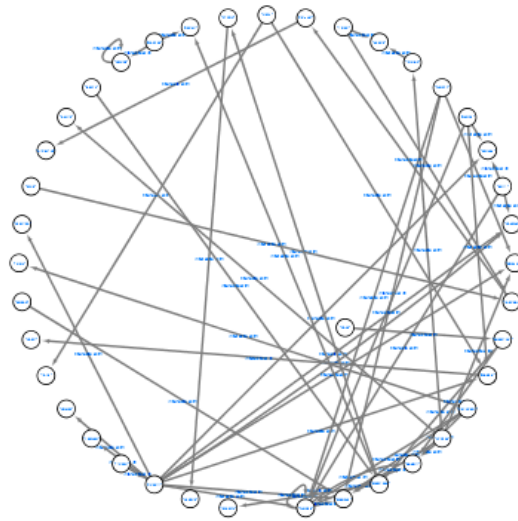


Figure 13 Prostate Cancer Regulatory Network.

Table 1 Dataset Profile

S. No.	Parameter	Value
1	Sample Count	24
2	Value Type	Transformed Count
3	Channel Count	1
4	Platform Organism	Homo sapiens
5	Platform	<i>In situ</i>
	Technology	<i>oligonucleotide</i>
6	Sample Type	RNA
7	Feature Count	54675
8	Dataset Platform	GPL570
9	Dataset identification	GDS4114
10	Series	GSE26910

Table 2 Intersecting transcripts in breast and prostate data.

S.No	ID_REF	Gene Symbol	Gene Title
1	1552509_a_	CD300LG	CD300 molecule-like family member g
2	203407_at	PPL	Periplakin
3	208891_at	DUSP6	dual specificity phosphatase 6
4	208892_s_a	DUSP6	dual specificity phosphatase 6
5	209426_s_a	AMACR	///alpha-
t		C1QTNF3	methylacyl-CoA racemase /// C1q and tumor necrosis factor related protein 3
6	209793_at	GRIA1	glutamate receptor,

ionotropic,

AMPA 1

7 210556_at NFATC3

nuclear factor of

activated T-cells,

cytoplasmic,

calcineurin-

dependent 3

Table 3 Breast cancer network visualization ready tabulation.

Source/From Gene	Target/To Gene	Correlation	P-value
PAX8	PAX8	-0.96158	6.18E-07
DDR1	AFG3L1P	-0.96217	5.72E-07
ZDHHC11	ALG10	0.965852	3.45E-07
C15orf40	PRSS33	0.957133	1.06E-06
TTC39C	TIRAP	-0.96351	4.79E-07
PXK	MSI2	-0.95376	1.54E-06
CORO6	FAM71A	-0.95738	1.03E-06
GIMAP1	GAPT	0.967306	2.78E-07
SPATA17	TSSK3	0.974818	7.64E-08
ENTHD1	CLEC12A	0.955864	1.22E-06
CENPBD1	C15orf27	0.970408	1.70E-07
WFDC2	CALML6	0.962181	5.72E-07
EYA3	DEFB106A ///	0.953635	1.56E-06
	DEFB106B		
CCDC65	DEFB106A ///	0.970579	1.65E-07
	DEFB106B		
MFAP3	C10orf25	0.958966	8.55E-07
TMEM106A	ETV3	0.955541	1.27E-06
KLHL10	KLHL10	0.969253	2.06E-07
RFC2	TM2D3	-0.95251	1.76E-06
SLC39A13	TM2D3	-0.96026	7.30E-07
C19orf26	TM2D3	-0.95676	1.11E-06
PRSS33	ANKAR	0.956311	1.16E-06
SLC39A13	SCGB1C1	0.952224	1.81E-06

CCDC65	SCGB1C1	-0.96076	6.86E-07
FAM122C	SCGB1C1	0.972535	1.18E-07
ADAM32	RAPH1	-0.96305	5.10E-07
PLCD3	SMCR8	-0.96539	3.69E-07
ARMCX4	BNC1	-0.95534	1.30E-06
LACTB	MPP4	0.958924	8.59E-07
DDR1	LACE1	0.952876	1.69E-06
SLC39A13	LACE1	-0.95083	2.08E-06
BRF1	LACE1	0.950253	2.21E-06
ZNF485	LACE1	0.975724	6.37E-08
TTLL12	IDI2	0.951264	1.99E-06
ARMCX4	CYP11B1	0.972441	1.20E-07
RAX2	TAF8	-0.95835	9.20E-07
TMEM106A	BRSK1	-0.95993	7.61E-07
TTC39C	KCNE4	0.965008	3.90E-07
CILP2	KCNE4	-0.95681	1.10E-06
RAX2	KCNE4	-0.95326	1.62E-06
COBL	KCNE4	0.950202	2.22E-06
CCL5	HIPK1	0.952769	1.71E-06
C19orf26	MTBP	0.966859	2.98E-07
ACAP2	MTBP	0.95023	2.21E-06
C19orf26	TMEM74	0.970411	1.70E-07
ZNF485	TMEM74	-0.98251	1.25E-08
IDI2	TMEM74	-0.97575	6.34E-08
IDI2	C21orf67	0.950287	2.20E-06
TMEM74	C21orf67	-0.95852	9.02E-07
ZDHHC11	NLRP11	0.950934	2.06E-06
BRSK1	NLRP11	-0.9579	9.70E-07

CCDC11	ATP6V1C2	-0.95934	8.17E-07
BRF1	ATP6V1C2	-0.95839	9.16E-07
RAPH1	TAGAP	-0.95451	1.42E-06
TM2D3	PRSS36	0.977508	4.37E-08
KCNE4	ZDHHC15	0.966886	2.97E-07
ATP6V1E2	CDK15	0.958811	8.71E-07
PRSS33	CDK15	-0.96044	7.14E-07
GAPT	CDK15	0.956134	1.19E-06
IDI2	CDK15	0.950309	2.19E-06
ZSCAN20		0.956726	1.11E-06
TMEM74		0.952653	1.73E-06
WFDC9	RTP3	0.950728	2.10E-06
RFC2	MIPOL1	0.951615	1.93E-06
SPATA17	MIPOL1	0.953921	1.51E-06
MEGF11	MIPOL1	0.954858	1.37E-06
PRSS33	MYO3B	-0.95153	1.94E-06
MEGF11	TRIML2	0.951737	1.90E-06
C8orf47	ABCC13	0.950995	2.05E-06
C21orf67	IL12RB1	-0.97209	1.27E-07
CILP2	GTF2A1L	0.956402	1.15E-06
PRSS33	GTF2A1L	0.953581	1.57E-06
MEGF11	GTF2A1L	0.952955	1.68E-06
TIGD4	GTF2A1L	-0.97121	1.48E-07
GTF2A1L	GTF2A1L	0.972378	1.21E-07
ACAP2	PXT1	0.951524	1.94E-06
LETM2	PXT1	0.950738	2.10E-06
PRUNE2	CDC42SE2	-0.97319	1.04E-07
NCRNA00204	CDC42SE2	0.95647	1.14E-06

ZNF485	RFWD2	-0.95757	1.01E-06
TMEM74	RFWD2	0.969177	2.08E-07
DDR1	STX6	-0.96476	4.03E-07
KCNE4	STX6	-0.95893	8.59E-07
RFC2	ANLN	-0.96962	1.94E-07
KLHL10	ANLN	0.954726	1.39E-06
FAM71A	TMEM163	-0.95918	8.33E-07
ZSCAN20	HERPUD2	-0.95071	2.11E-06
KLK8	JMJD6	0.953033	1.66E-06
CATSPER1	MAP3K6	0.983465	9.47E-09
CILP2	PTPN11	0.961084	6.58E-07
PRSS33	PTPN11	0.95872	8.81E-07
TTLL10	PTPN11	0.951342	1.98E-06
MEGF11	PTPN11	0.957441	1.02E-06
MIPOL1	PTPN11	0.953438	1.59E-06
ABCC13	PTPN11	0.975157	7.15E-08
RDH10	KLHDC7B	0.960543	7.05E-07
CCDC65	PHC3	0.976601	5.31E-08
GIMAP1	TNFRSF10A	0.976163	5.82E-08
FLJ30901	TNFRSF10A	0.971343	1.45E-07
KLHL10	RFFL	0.950583	2.14E-06
NEXN	RFFL	0.955226	1.31E-06
NEXN	HPS4	0.97026	1.74E-07
HPS4	HPS4	0.970347	1.72E-07
CCL5	UHMK1	0.972703	1.14E-07
HIPK1	UHMK1	0.955461	1.28E-06
CLEC12A	TXNDC2	0.96067	6.94E-07
CCL5	C5orf22	0.956595	1.13E-06

IDI2	PCDHGB7	0.950937	2.06E-06
GPBAR1	ERC1	0.958682	8.85E-07
KLHL10	FLCN	0.950509	2.15E-06
SPRR4	FLCN	0.970464	1.68E-07
GUCA1A	SLC9A7	0.956361	1.16E-06
PRUNE2	DIRC1	0.967442	2.73E-07
NCRNA00204	DNAJB7	0.953422	1.60E-06
ETV3	CASC5	0.950507	2.15E-06
DDR1	C20orf152	-0.95261	1.74E-06
CORO6	C20orf152	0.957301	1.04E-06
GAPT	C20orf152	-0.96149	6.25E-07
IDI2	C20orf152	-0.96458	4.14E-07
TMEM74	C20orf152	0.966728	3.04E-07
TMEM106A	CASKIN1	0.985543	4.85E-09
BSND	CACNA2D4	-0.95795	9.65E-07
ERC1	MGC16703	0.954553	1.41E-06
ERC1	CARD16 ///	0.954911	1.36E-06
	CASP1		
IDI2	DUSP19	0.951786	1.89E-06
ZNF570	DUSP19	0.9683	2.39E-07
LACE1	CYB5D1	0.976222	5.75E-08
C19orf26	CREG2	0.967519	2.70E-07
ETV3	CREG2	0.957063	1.07E-06
TM2D3	RXFP1	0.969355	2.02E-07
KLHL10	SPEF2	-0.95268	1.73E-06
DEFB106A ///	DTD1	0.951094	2.03E-06
DEFB106B			
ABCC13	DTD1	0.957523	1.01E-06

VPS18	CASC4	-0.96446	4.21E-07
CACNG5	CASC4	0.953388	1.60E-06
FAM122C	FGF1	0.953203	1.63E-06
ALG10	ARPP21	0.961554	6.20E-07
BRF1	ARPP21	0.958093	9.49E-07
ZNF485	ARPP21	0.964801	4.01E-07
IDI2	ARPP21	0.984576	6.70E-09
TMEM74	ARPP21	-0.96692	2.95E-07
CLDN19	ARPP21	-0.95236	1.78E-06
ZNF570	ARPP21	0.956583	1.13E-06
C21orf29	ARPP21	0.956154	1.19E-06
HPS4	ARPP21	0.95544	1.28E-06
CASKIN1	ARPP21	-0.95654	1.14E-06
NEDD1	ADAMTS17	0.965781	3.49E-07
GIMAP1	ADAMTS17	-0.95222	1.81E-06
BSND	ADAMTS17	-0.96552	3.62E-07
ARL11	ADAMTS17	0.964328	4.28E-07
ADAMTS17	ADAMTS17	0.953315	1.61E-06
EPHB3	DHH	-0.95225	1.80E-06
EPHB3	KLHDC1	-0.96504	3.88E-07
SERPINA12	RICTOR	0.957731	9.90E-07
ARPP21	RICTOR	-0.95717	1.06E-06
C8orf47	PCDHGA4	-0.96422	4.35E-07
NCRNA00161	PCDHGA4	-0.96077	6.85E-07
C21orf67	NETO1	-0.95729	1.04E-06
C5orf22	NETO1	0.9569	1.09E-06
LACTB	ST7L	-0.95159	1.93E-06

Table 4 Prostate cancer network visualization ready tabulation.

Source/From Gene	Target/To Gene	Correlation	P-value
BEST4	TMEM106A	-0.96094	6.70E-07
CYP2A6	ALG10	0.959298	8.22E-07
C15orf40	C15orf40	0.982069	1.42E-08
TTC39C	CCDC11	-0.95256	1.75E-06
CYP2A6	TRIOBP	-0.95726	1.05E-06
C15orf40	CRYZL1	-0.95591	1.22E-06
TRIOBP	LEAP2	0.951457	1.96E-06
TIRAP	LEAP2	0.970564	1.66E-07
FAM122C	SCIN	-0.97165	1.38E-07
TIRAP	FAM18B2	-0.95748	1.02E-06
MSI2	FAM18B2	-0.97047	1.68E-07
PRR22	FAM71A	0.957059	1.07E-06
PXK	FAM71A	0.957061	1.07E-06
CCDC65	FAM71A	-0.95069	2.11E-06
SCIN	GAPT	0.960974	6.68E-07
DDR1	C8orf47	0.953643	1.56E-06
TIMD4	C1orf65	-0.95672	1.11E-06
CCDC11	CLEC12A	-0.95066	2.12E-06
CCDC11	CALML6	-0.96175	6.05E-07
CCDC11	CALML6	-0.97601	6.01E-08
BRF1	CALML6	0.957768	9.85E-07
GIMAP1	DEFB106A ///	0.960782	6.84E-07
	DEFB106B		
CCDC65	DEFB106A ///	0.951651	1.92E-06

	DEFB106B		
RDH10	DEFB106A ///	0.968938	2.16E-07
	DEFB106B		
VPS18	WFDC9	-0.95189	1.87E-06
CCDC11	ZNF485	0.968029	2.49E-07
BRF1	ZNF485	-0.97176	1.35E-07
NLRP5	ZNF485	-0.95257	1.74E-06
CCDC11	CDH23	-0.96333	4.91E-07
CYP2A6	DNAJC5G	0.960958	6.69E-07
CCDC11	DNAJC5G	-0.95102	2.04E-06
WDR17	DNAJC5G	-0.95322	1.63E-06
CCDC11	WDR88	-0.96017	7.38E-07
CATSPER1	WDR88	-0.97324	1.03E-07
BRF1	WDR88	0.954811	1.38E-06
NLRP5	WDR88	0.963964	4.50E-07
WDR17	WDR88	-0.96062	6.98E-07
WDR88	WDR88	0.978602	3.41E-08
CYP2A6	MBD3L2	-0.96108	6.59E-07
CCDC11	MBD3L2	0.971062	1.52E-07
CATSPER1	MBD3L2	0.981525	1.64E-08
WDR17	MBD3L2	0.965176	3.80E-07
ANKAR	MBD3L2	-0.96971	1.91E-07
WDR88	MBD3L2	-0.95909	8.42E-07
WDR88	MBD3L2	-0.96537	3.70E-07
PAX8	ADAMTSL1	0.950173	2.22E-06
TMEM106A	ADAMTSL1	-0.96193	5.91E-07
ODF4	MBD3L1	-0.95259	1.74E-06
CCDC11	MBD3L1	0.988363	1.65E-09

NLRP5	MBD3L1	-0.95047	2.16E-06
CCDC11	FAM46D	0.950902	2.07E-06
ADAM32	SERPINB11	-0.96303	5.11E-07
NEDD1	DSCR10	0.957306	1.04E-06
CYP2A6	ABHD11	-0.95662	1.13E-06
CATSPER1	ABHD11	0.965901	3.43E-07
WDR88	ABHD11	-0.96721	2.82E-07
WDR88	ABHD11	-0.96596	3.40E-07
ADAMTSL1	ABHD11	0.953065	1.66E-06
MBD3L1	GAMT	-0.95202	1.85E-06
TMEM106A	PTPRC	-0.9591	8.41E-07
TMEM106A	RAPH1	-0.95232	1.79E-06
MAN1A2	RAPH1	0.993202	1.13E-10
MAN1A2	SMCR8	0.993981	6.16E-11
SMCR8	SMCR8	0.998425	7.62E-14
PDE7A	LACTB	0.95717	1.06E-06
BNC1	BNC1	0.951039	2.04E-06
ODF4	TAF8	0.951929	1.86E-06
KLK8	SLAMF6	-0.96746	2.72E-07
SCARB1	ZSCAN20	-0.95303	1.66E-06
C4orf33	RHBDL2	-0.95553	1.27E-06
SERPINB11	RHBDL2	0.951318	1.98E-06
ARMCX4	CRB2	-0.95373	1.55E-06
FAM122C	WBP2NL	-0.97404	8.89E-08
TMEM106A	TMEM74	-0.95533	1.30E-06
CATSPER1	TIGD4	0.960038	7.50E-07
FAM18B2	C21orf67	0.989391	1.04E-09
FAM71A	NLRP11	0.966506	3.14E-07

CLEC4F	NLRP11	0.965087	3.85E-07
C21orf67	ATP6V1C2	0.957045	1.07E-06
PTPRC	CLDN19	-0.96369	4.68E-07
TIMD4	KIF6	0.952698	1.72E-06
DDR1	TAGAP	0.95111	2.03E-06
PRR22	TAGAP	0.956625	1.12E-06
HIPK1	TAGAP	0.96528	3.75E-07
KLHL10	LETM2	-0.96726	2.80E-07
ESX1	BSND	0.965532	3.62E-07

Stromal Data Analysis: R Script file.

```
# Installing GEOquery
```

```
source("http://www.bioconductor.org/biocLite.R")
```

```
biocLite("GEOquery")
```

```
# Loading GEO file with GEOquery
```

```
library(Biobase)
```

```
library(GEOquery)
```

```
#Download GPL file, put it in the current directory, and load it:
```

```
gpl570 <- getGEO('GPL570', destdir=".")
```

```
#Or, open an existing GPL file:
```

```
gpl570 <- getGEO(filename='GPL570.soft')
```

```
# Handpicked description (three columns: ID, Gene Symbol, Gene Title).
```

```
Table(gpl570) [c("ID", "Gene Symbol", "Gene Title")]
```

```
IDs <- attr(dataTable(gpl570), "table")[, c("ID", "Gene Symbol", "Gene Title")]
```

```
# Extract the expression values from the dataset
```

```
# line 64 contains field names
```

```
DS_Main <- read.table("GSE26910_series_matrix.txt.gz", skip = 63, header = TRUE, sep = "\t", fill = TRUE)
```

```
# Remove the last line from the matrix that says "!series_matrix_table_end"
```

```
DS_Main <- DS_Main[-54676, ]
```

```
# Merging the annotation information to the expression values matrix and rejecting null entries.
```

```
names(IDs)[1] <- "ID_REF"
```

```
DS <- merge(IDs,DS_Main, by = "ID_REF")
```

```
DS[DS == ""] <- NA
```

```
DS <- na.omit(DS)
```

```
# Reordering of respective breast cancer and prostate cancer datasets.
```

```
# Prostate Normal [1:6], Prostate Tumor [7:12], Breast Normal [13:18], Breast Tumor [19:24]
```

```
WorkDS <- DS [c(4,6,8,10,12,14, 5,7,9,11,13,15, 16,18,20,22,24,26, 17,19,21,23,25,27)]
```

```
# RowMeans calculation
```

```
ProstateNormalMean <- rowMeans(log2(WorkDS[,1:6]))
```

```
ProstateTumorMean <- rowMeans(log2(WorkDS[,7:12]))
```

```
BreastNormalMean <- rowMeans(log2(WorkDS[,13:18]))
```

```
BreastTumorMean <- rowMeans(log2(WorkDS[,19:24]))
```

```
# MA-Plot
```

```
par(mfrow=c(1,2))
```

```
ProstateMean <- rowMeans(log2(WorkDS[, 1:12]))
```

```
BreastMean <- rowMeans(log2(WorkDS[, 13:24]))
```

```
plot(ProstateMean, ProstateTumorMean-ProstateNormalMean, main="MA Plot for Prostate data", pch=16,
```

```
cex=0.35)
```

```
hold()
```

```
plot(BreastMean, BreastTumorMean-BreastNormalMean, main="MA Plot for Breast data", pch=16, cex=0.35)
```

```
# Rough draft of extreme probes
```

```
DS[which.min(BreastTumorMean-BreastNormalMean), ] ### most negatively expressed breast gene
```

```
DS[which.min(ProstateTumorMean-ProstateNormalMean), ] ### most negatively expressed prostate gene
```

```
DS[which.max(ProstateTumorMean-ProstateNormalMean), ] ### most positively expressed prostate gene
```

```
DS[which.max(BreastTumorMean-BreastNormalMean), ] ### most positively expressed breast gene
```

```
# Standard Deviation calculation for t-test
```

```
install.packages(genefilter)
```

```
library(genefilter)
```

```
ProstateNormalSD <- rowSds(log2(WorkDS[,1:6]))
```

```
ProstateTumorSD <- rowSds(log2(WorkDS[,7:12]))
```

```
BreastNormalSD <- rowSds(log2(WorkDS[,13:18]))
```

```
BreastTumorSD <- rowSds(log2(WorkDS[,19:24]))
```

```
# t-test calculation and histogram plot
```

```
par(mfrow=c(1,2))
```

```
Prostate_ttest <- (ProstateTumorMean-ProstateNormalMean)/sqrt(ProstateTumorSD^2/6 +
```

```
ProstateNormalSD^2/6)
```

```
hist(Prostate_ttest,nclass=100)
```

```
hold()
```

```
Breast_ttest <- (BreastTumorMean-BreastNormalMean)/sqrt(BreastTumorSD^2/6 + BreastNormalSD^2/6)
```

```
hist(Breast_ttest, nclass=100)
```

```
# p-value calculation and histogram plot
```

```
Prostate_pval <- 2*(1-pt(abs(Prostate_ttest),5))
```

```
Breast_pval <- 2*(1-pt(abs(Breast_ttest),5))
```

```
par(mfrow=c(1,2))
```

```
hist(Prostate_pval, nclass=100)
```

```
hold()
```

```
hist(Breast_pval, nclass = 100)
```

```
# volcano Plot
```

```
par(mfrow=c(1,2))
```

```
plot(ProstateTumorMean-ProstateNormalMean, -log10(Prostate_pval), main = "Volcano Plot@Prostate tissue",  
xlab = "Sample Mean Difference", ylab = "-log10(p value)", pch=16, cex=0.35)
```

```
hold()
```

```
plot(BreastTumorMean-BreastNormalMean, -log10(Breast_pval), main = "Volcano Plot@Breast tissue", xlab=  
"Sample Mean Difference", ylab = "-log10(p value)", pch=16, cex=0.35)
```

```
# Boxplots for the normal data and its log transformed version.(Log2 transformation applied)
```

```
par(mfrow = c(1, 2))
```

```
boxplot(WorkDS, col = c(2,3,2,3,2,3,2,3,2,3,2,3), main = "Expression values pre-normalization",  
xlab = "Slides", ylab = "Expression", las = 2, cex.axis = 0.7)
```

```
hold()
```

```
boxplot(log2(WorkDS), col = c(2,3,2,3,2,3,2,3,2,3,2,3), main = "Expression values post-log-transformation",  
xlab = "Slides", ylab = "Expression", las = 2, cex.axis = 0.7)
```

```
abline(0, 0, col = "black")
```

```
# Check Normality
```

```
par(mfrow=c(1,2))
```

```
qqnorm(Prostate_ttest, main = "QQ Plot@Prostate Data")
```

```
qqline(Prostate_ttest)
```

```
hold()
```

```
qqnorm(Breast_ttest, main = "QQ Plot@Breast Data")
```

```
qqline(Breast_ttest)
```

```
# Elucidating genes with particular p-values.
```

```
for (i in c(0.01, 0.05, 0.001, 1e-04, 1e-05, 1e-06, 1e-07))
```

```
  print(paste("genes with p-values smaller than",i, length(which(Prostate_pval < i))))
```

```
for (i in c(0.01, 0.05, 0.001, 1e-04, 1e-05, 1e-06, 1e-07))
```

```
  print(paste("genes with p-values smaller than",i, length(which(Breast_pval < i))))
```

```
# Plot heatmap of differentially expressed genes: Genes are differentially expressed if its p-value is under a given threshold, which must be smaller than the usual 0.05 or 0.01 due to multiplicity of tests
```

```
BreastDEGenes <- data.frame(which(Breast_pval < 0.01))
```

```
ProstateDEGenes <- data.frame(which(Prostate_pval < 0.01))
```

```
ProstateDEGenesData <- ProstateDEGenes[,1]

BreastDEGenesData <- BreastDEGenes[,1]

ProstateData <- as.matrix(WorkDS[ProstateDEGenesData, 1:12])

heatmap(ProstateData, col = topo.colors(100), cexRow = 0.5)

BreastData <- as.matrix(WorkDS[BreastDEGenesData, 13:24])

heatmap(BreastData, col = topo.colors(100), cexRow = 0.5)

# List of differentially expressed genes.

#Breast Data

BDEG <- matrix(nrow = nrow(BreastDEGenes), ncol = 1)

for(i in 1:nrow(BreastDEGenes)) BDEG[i,]<- paste(DS[BreastDEGenes[i,], "ID_REF"])

BDEG <- as.data.frame(BDEG)

names(BDEG)[1] <- "ID_REF"

FinalBDEG <- merge(BDEG,DS)

BDEG <- merge(BDEG, IDs, by = 'ID_REF')

view(BDEG)

#Prostate Data

PDEG <- matrix(nrow = nrow(ProstateDEGenes),ncol = 1)

for(i in 1:nrow(ProstateDEGenes)) PDEG[i,] <- paste(DS[ProstateDEGenes[i,], "ID_REF"])

PDEG <- as.data.frame(PDEG)

names(PDEG)[1] <- "ID_REF"

FinalPDEG <- merge(PDEG,DS)

PDEG <- merge(PDEG, IDs, by = 'ID_REF')
```



```
view(PDEG)
```

```
##Intersecting transcripts in breast and prostate cancer types as marked in the dataset
```

```
BDEG$match <- match(BDEG$location, PDEG$location, nomatch=0)
```

```
# Reordering of respective breast cancer and prostate cancer datasets.
```

```
# Prostate Normal [1:6], Prostate Tumor [7:12], Breast Normal [13:18], Breast Tumor [19:24]
```

```
FinalPDEG <- FinalPDEG [c(4,6,8,10,12,14, 5,7,9,11,13,15, 16,18,20,22,24,26, 17,19,21,23,25,27)]
```

```
WorkFinalPDEG <- FinalPDEG[1:12]
```

```
FinalBDEG <- FinalBDEG [c(4,6,8,10,12,14, 5,7,9,11,13,15, 16,18,20,22,24,26, 17,19,21,23,25,27)]
```

```
WorkFinalBDEG <- FinalBDEG[13:24]
```

```
##Prostate data regrerssion analysis(linear model)
```

```
par(mfrow=c(1,6))
```

```
plot(log2(WorkFinalPDEG$GSM662756),log2(WorkFinalPDEG$GSM662757), pch = 16, cex = 1.3, col =  
c("blue", "red"))
```

```
abline(lm(log2(WorkFinalPDEG$GSM662756) ~ log2(WorkFinalPDEG$GSM662757)), col= 1)
```

```
plot(log2(WorkFinalPDEG$GSM662758),log2(WorkFinalPDEG$GSM662759), pch = 16, cex = 1.3, col =  
c("blue", "red"))
```

```
abline(lm(log2(WorkFinalPDEG$GSM662758) ~ log2(WorkFinalPDEG$GSM662759)), col= 1)
```

```
plot(log2(WorkFinalPDEG$GSM662760),log2(WorkFinalPDEG$GSM662761), pch = 16, cex = 1.3, col =  
c("blue","red"))  
abline(lm(log2(WorkFinalPDEG$GSM662760) ~ log2(WorkFinalPDEG$GSM662761)), col= 1)  
plot(log2(WorkFinalPDEG$GSM662762),log2(WorkFinalPDEG$GSM662763), pch = 16, cex = 1.3, col =  
c("blue","red"))  
abline(lm(log2(WorkFinalPDEG$GSM662762) ~ log2(WorkFinalPDEG$GSM662763)), col= 1)  
plot(log2(WorkFinalPDEG$GSM662764),log2(WorkFinalPDEG$GSM662765), pch = 16, cex = 1.3, col =  
c("blue","red"))  
abline(lm(log2(WorkFinalPDEG$GSM662764) ~ log2(WorkFinalPDEG$GSM662765)), col= 1)  
plot(log2(WorkFinalPDEG$GSM662766),log2(WorkFinalPDEG$GSM662767), pch = 16, cex = 1.3, col =  
c("blue","red"))  
abline(lm(log2(WorkFinalPDEG$GSM662766) ~ log2(WorkFinalPDEG$GSM662767)), col= 1)
```

##Gene Set Enrichment Analysis

```
library(genefilter)
```

```
library(GSEABase)
```

```
Breast_GSEA <- GeneSetCollection(WorkFinalBDEG, setType = KEGGCollection())
```

```
Prostate_GSEA <- GeneSetCollection(WorkFinalPDEG, setType = KEGGCollection())
```

##Correlation Analysis

##Breast

```
WorkFinalBDEG <- read.csv("WorkFinalBDEG_GSEAFiltered.csv") ## Import filtered annotation file from  
MeV.
```

```
btemp <- WorkFinalBDEG
```

```
btemp$ID_REF <- NULL

btemp <- log2(btemp)

pairs(btemp)

BreastCorrelationMatrix <- cor(t(as.matrix(btemp)))

BreastCorMat <- as.data.frame(BreastCorrelationMatrix)

rownames(BreastCorMat) <- WorkFinalBDEG$ID_REF

colnames(BreastCorMat) <- WorkFinalBDEG$ID_REF

##Prostate

WorkFinalPDEG <- read.csv("WorkFinalPDEG_GSEAFiltered.csv") ## Import filtered annotation file from
MeV.

ptemp <- WorkFinalPDEG

ptemp$ID_REF <- NULL

ptemp <- log2(ptemp)

pairs(ptemp)

ProstateCorrelationMatrix <- cor(t(as.matrix(ptemp)))

ProCorMat<- as.data.frame(ProstateCorrelationMatrix)

rownames(ProCorMat) <- WorkFinalPDEG$ID_REF

colnames(ProCorMat)<- WorkFinalPDEG$ID_REF

### Feature Selection: Clustering of robustly entwined genes.

install.packages("gplots")

install.packages("Hmisc")

library(Hmisc)

library(gplots)
```

```
heatmap.2(ProstateCorrelationMatrix, main="Hierarchical Cluster",
dendrogram="column",trace="none",col=greenred(10))
heatmap.2(1-abs(ProstateCorrelationMatrix), distfun=as.dist, trace="none")
heatmap.2(BreastCorrelationMatrix, main="Hierarchical Cluster",
dendrogram="column",trace="none",col=greenred(10))
heatmap.2(1-abs(BreastCorrelationMatrix), distfun=as.dist, trace="none")
```

```
##Prostate Data
```

```
library(caret)
```

```
HighlyCorrelated <- findCorrelation(ProstateCorrelationMatrix, cutoff = 0.95, verbose = TRUE, names =
FALSE)
```

```
print(HighlyCorrelated)
```

```
WorkFinalPDEG[HighlyCorrelated,1]
```

```
IDs[WorkFinalPDEG[HighlyCorrelated,1],c(2,3)]
```

```
for(i in 2:nrow(BreastCorMat))
```

```
{
```

```
  for(j in 1:ncol(BreastCorMat)-1)
```

```
  {
```

```
    if(i>j)
```

```
    {
```

```
      out <- c (rownames(BreastCorMat[i,]), colnames(BreastCorMat[j]), BreastCorMat[i,j])
```

```
      write.table(out, file="output.txt", append=TRUE, sep= " ")
```

```
    }
```

```
  else
```

```
    break
  }
}

#Network Ready Matrix Format (Function) // Credit: http://www.sthda.com

flattenCorrMatrix <- function(cmat, pmat) {
  ut <- upper.tri(cmat)
  data.frame(
    row = rownames(cmat)[row(cmat)[ut]],
    column = rownames(cmat)[col(cmat)[ut]],
    cor = (cmat)[ut],
    p = pmat[ut]
  )
}

library(Hmisc)

btemp <- as.matrix(btemp)
rownames(btemp) <- WorkFinalBDEG$ID_REF
BreastNet <- rcorr(t(btemp))
BreastNetworkInputMatrix <- flattenCorrMatrix(BreastNet$r, BreastNet$p)

#lets map the gene names to row and column entries
BreastNetworkInputMatrix$row <- IDs[WorkFinalBDEG[BreastNetworkInputMatrix$row,1],2]
BreastNetworkInputMatrix$column <- IDs[WorkFinalBDEG[BreastNetworkInputMatrix$column,1],2]
```

```
p-temp <- as.matrix(p-temp)

rownames(p-temp) <- WorkFinalPDEG$ID_REF

ProstateNet <- rcorr(t(p-temp))

ProstateNetworkInputMatrix <- flattenCorrMatrix(ProstateNet$r, ProstateNet$P)

ProstateNetworkInputMatrix$row <- IDs[WorkFinalPDEG[ProstateNetworkInputMatrix$row,1],2]

ProstateNetworkInputMatrix$column <- IDs[WorkFinalPDEG[ProstateNetworkInputMatrix$column,1],2]

symnum(BreastCorrelationMatrix)

symnum(ProstateCorrelationMatrix)

install.packages("corrplot")

library(corrplot)

corrplot(BreastCorrelationMatrix, type="upper", order="hclust", tl.col="black", tl.srt=45)

corrplot(ProstateCorrelationMatrix, type="upper", order="hclust", tl.col="black", tl.srt=45)

install.packages("PerformanceAnalytics")

library(PerformanceAnalytics)

chart.Correlation(BreastCorrelationMatrix, histogram= TRUE, pch= 19)

chart.Correlation(ProstateCorrelationMatrix, histogram= TRUE, pch= 19)

col <- colorRampPalette(c("blue", "white", "red"))(20)

heatmap(x = BreastCorrelationMatrix, col = col, symm = TRUE)

heatmap(x = ProstateCorrelationMatrix, col = col, symm = TRUE)

## Optimize network ready correlation and p-values matrix

## top candidates which manifest low p-value and high correlation.
```

```
BreastFinal <- BreastNetworkInputMatrix[which(abs(BreastNetworkInputMatrix$cor) > 0.95 |  
BreastNetworkInputMatrix$p < 0.0000001), c(1,2,3,4)]  
ProstateFinal <- ProstateNetworkInputMatrix[which(abs(ProstateNetworkInputMatrix$cor) > 0.95 |  
ProstateNetworkInputMatrix$p < 0.0000001), c(1,2,3,4)]  
write.csv(BreastFinal, "BreastFinalTest.csv")  
write.csv(ProstateFinal, "ProstateFinalTest.csv")  
  
##Intersecting transcripts in breast and prostate cancer types as marked in the dataset  
  
BDEG$match <- match(BDEG$location, PDEG$location, nomatch=0)  
  
# Reordering of respective breast cancer and prostate cancer datasets.  
# Prostate Normal [1:6], Prostate Tumor [7:12], Breast Normal [13:18], Breast Tumor [19:24]  
  
FinalPDEG <- FinalPDEG [c(4,6,8,10,12,14, 5,7,9,11,13,15, 16,18,20,22,24,26, 17,19,21,23,25,27)]  
WorkFinalPDEG <- FinalPDEG[1:12]  
  
FinalBDEG <- FinalBDEG [c(4,6,8,10,12,14, 5,7,9,11,13,15, 16,18,20,22,24,26, 17,19,21,23,25,27)]  
WorkFinalBDEG <- FinalBDEG[13:24]  
  
##Prostate data regrerssion analysis(linear model)  
  
par(mfrow=c(1,6))  
plot(log2(WorkFinalPDEG$GSM662756),log2(WorkFinalPDEG$GSM662757), pch = 16, cex = 1.3, col =  
c("blue","red"))
```

```
abline(lm(log2(WorkFinalPDEG$GSM662756) ~ log2(WorkFinalPDEG$GSM662757)), col= 1)
plot(log2(WorkFinalPDEG$GSM662758),log2(WorkFinalPDEG$GSM662759), pch = 16, cex = 1.3, col =
c("blue","red"))
abline(lm(log2(WorkFinalPDEG$GSM662758) ~ log2(WorkFinalPDEG$GSM662759)), col= 1)
plot(log2(WorkFinalPDEG$GSM662760),log2(WorkFinalPDEG$GSM662761), pch = 16, cex = 1.3, col =
c("blue","red"))
abline(lm(log2(WorkFinalPDEG$GSM662760) ~ log2(WorkFinalPDEG$GSM662761)), col= 1)
plot(log2(WorkFinalPDEG$GSM662762),log2(WorkFinalPDEG$GSM662763), pch = 16, cex = 1.3, col =
c("blue","red"))
abline(lm(log2(WorkFinalPDEG$GSM662762) ~ log2(WorkFinalPDEG$GSM662763)), col= 1)
plot(log2(WorkFinalPDEG$GSM662764),log2(WorkFinalPDEG$GSM662765), pch = 16, cex = 1.3, col =
c("blue","red"))
abline(lm(log2(WorkFinalPDEG$GSM662764) ~ log2(WorkFinalPDEG$GSM662765)), col= 1)
plot(log2(WorkFinalPDEG$GSM662766),log2(WorkFinalPDEG$GSM662767), pch = 16, cex = 1.3, col =
c("blue","red"))
abline(lm(log2(WorkFinalPDEG$GSM662766) ~ log2(WorkFinalPDEG$GSM662767)), col= 1)
```

##Gene Set Enrichment Analysis

```
library(genefilter)
```

```
library(GSEABase)
```

```
Breast_GSEA <- GeneSetCollection(WorkFinalBDEG, setType = KEGGCollection())
```

```
Prostate_GSEA <- GeneSetCollection(WorkFinalPDEG, setType = KEGGCollection())
```

Support Vector Machine Implementation

```
## Prostate
```



```
install.packages("e1071")

library(e1071)

temp1 <- WorkFinalPDEG

temp1$ID_REF <- NULL

temp1 <- log2(temp1)

temp1 <- t(temp1)

ClassLabels1 <- c(rep(1,6),rep(-1,6))

DataFrame1 <- data.frame(Gene=temp1,ClassLabels=as.factor(ClassLabels1))

SVMModel1 <- svm(ClassLabels1~., data=DataFrame1, kernel="linear", cost=10, scale = FALSE)

GeneWeights1<-t(SVMModel1$coefs)%*%SVMModel1$SV

sort.list(GeneWeights1) ## Genes 212 and 129 have highest and second highest weights, respectively.

plot(SVMModel1,DataFrame1, Gene.212 ~ Gene.129)

##Breast

temp2 <- WorkFinalBDEG

temp2$ID_REF <- NULL

temp2 <- log2(temp2)

temp2 <- t(temp2)

ClassLabels2<- c(rep(1,6),rep(-1,6))

DataFrame2 <- data.frame(Gene=temp2,ClassLabels=as.factor(ClassLabels2))

SVMModel2 <- svm(ClassLabels2~., data=DataFrame2, kernel="linear", cost=10, scale = FALSE)

GeneWeights2<-t(SVMModel2$coefs)%*%SVMModel2$SV

sort.list(GeneWeights2) ## Genes 346 and 133 have highest and second highest weights, respectively.

plot(SVMModel2,DataFrame2, Gene.346 ~ Gene.133)
```

```
install.packages("kernlab")
```

```
library(kernlab)
```

```
x <- as.matrix(temp1)
```

```
y <- matrix(c(rep(1,6),rep(-1,6)))
```

```
svp <- ksvm(x,y,type="C-svc", prob.model= TRUE)
```

```
predict (svp, x, type= "probabilities")
```