

Mining PubMed for biomarker-disease associations to guide discovery

Biomedical knowledge is growing exponentially; however, meta-knowledge around the data is often lacking. PubMed is a database comprising more than 21 million citations for biomedical literature from MEDLINE and additional life science journals dating back to the 1950s. To explore the use and frequency of biomarkers across human disease, we mined PubMed for biomarker-disease associations. We then ranked the top 100 linked diseases by relevance and mapped them to medical subject headings (MeSH) and, subsequently, to the Disease Ontology. To identify biomarkers for each disease, we queried Covance BioPathways, an online data resource that maps commercial biomarker assays to biological and disease pathways. We then integrated pathways-based information to describe both known and potential biomarkers as well as disease-associated genes/proteins for select diseases. This approach identifies therapeutic areas with candidate or validated biomarkers, and highlights those areas where a paucity of biomarkers exists.

Walter J. Jessen, Katherine T. Landschulz, Thomas G. Turi and Rachel Y. Reams
Covance Biomarker Center of Excellence, Discovery and Translational Services, Greenfield, Indiana

Biomedical knowledge is growing exponentially; however, meta-knowledge around the data is often lacking. PubMed is a database comprising more than 21 million citations for biomedical literature from MEDLINE and additional life science journals dating back to the 1950s. To explore the use and frequency of biomarkers across human disease, we mined PubMed for biomarker-disease associations. We then ranked the top 100 linked diseases by relevance and mapped them to medical subject headings (MeSH) and, subsequently, to the Disease Ontology. To identify biomarkers for each disease, we queried Covance BioPathways, an online data resource that maps commercial biomarker assays to biological and disease pathways. We then integrated pathways-based information to describe both known and potential biomarkers as well as disease-associated genes/proteins for select diseases. This approach identifies therapeutic areas with candidate or validated biomarkers, and highlights those areas where a paucity of biomarkers exists.

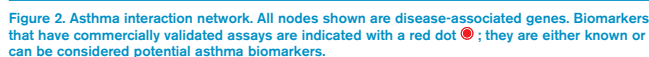
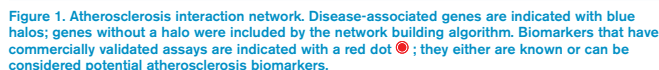
In June 2011, we mined PubMed for term ("biomarker")-disease associations and identified a total of 1,181 disease associations (Table 1). We then curated the top 100 disease associations from the list, mapping each result to both medical subject (MeSH) ID and Disease Ontology ID (DOID), and then subsequently queried the GeneGo diseases ontology for associated biomarkers (Table 2). Of 100 results, 62 map to both MeSH ID and DOID and are shown below.

Table 1. A representative list of term ("biomarker")-disease associations mined from PubMed in June 2011. The top 100 disease associations were ranked by Z Score. The Z-score indicates the number of standard deviations that the relevancy score is above the mean; larger Z-scores denote stronger associations. The top 100 data set is available under the Open Data Commons Attribution License at <http://BiomarkerCommons.org>.

Table 2. The curated list of disease associations mined from PubMed and organized by high-level Disease Ontology. Each specific disease association has a unique MeSH ID, DOID and number of associated genes as defined in the GeneGo MetaCore knowledgebase.

Text mining was performed using PolySearch, a web-based text mining system for extracting relationships between human diseases, genes, mutations, drugs and metabolites [Cheng et al., 2008]. The MeSH Browser (2012 MeSH) was used to map disease associations to MeSH IDs. Once MeSH IDs were assigned, the Disease Ontology was used to map DOIDs [Schriml et al., 2011]. Interaction networks were constructed in GeneGo MetaCore [Ekins et al., 2007] using the Auto expand algorithm, which gradually expands sub-networks around every object from the seed object list based on interactions identified in the literature. At every step, preference is given to objects with more connectivity to the initial object, and expansion halts when the sub-networks intersect, or when the overall network size reaches a predefined limit. Genes/proteins for which validated commercial assays exist were identified using Covance BioPathways at <http://www.Covance.com/BioPathways> and are indicated with a red dot ●. These genes/proteins can be considered potential biomarkers.

For illustrative purposes, we constructed an interaction network around disease-associated genes for two diseases—one with few associated genes (atherosclerosis) and one with many associated genes (asthma)—using a network building algorithm in GeneGo MetaCore. For each interaction network gene set, we then queried Covance BioPathways, a publicly accessible, web-based data source that integrates biological and disease pathway maps with validated Covance assays and antibody products, to identify commercially available biomarker assays.



Given the molecular interdependencies within a cell, a disease is rarely a consequence of a single gene abnormality but instead reflects the perturbation of a complex network of biological and signaling pathways. The approach described here describes the detection and ranking of human disease based on research/clinical activity surrounding biomarkers. It also enables the identification of therapeutic areas with candidate or validated biomarkers. The strategy takes an integrative approach to identify candidate disease biomarkers by combining disease-associated genes/proteins with commercially validated assays for known biomarkers. We first constructed a system-level model of disease that incorporates molecular interactions across biological and signaling pathways. We then identified each gene/protein in the model that has an existing commercially validated assay. This research offers an alternative, comprehensive view of key relationships and pathway perturbations that may identify biomarkers of disease emergence or progression.