

© 2015 Minji Kim

GENE PRIORITIZATION THROUGH  
HYBRID DISTANCE-SCORE RANK AGGREGATION

BY

MINJI KIM

THESIS

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Electrical and Computer Engineering  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2015

Urbana, Illinois

Adviser:

Associate Professor Olgica Milenkovic

# ABSTRACT

This thesis is concerned with developing novel rank aggregation methods for gene prioritization. Gene prioritization refers to a family of computational techniques for inferring disease genes through a set of training genes and carefully chosen similarity criteria. Test genes are scored based on their average similarity to the training set, and the rankings of genes under various similarity criteria are aggregated via statistical methods. The contributions of our work are threefold: a) First, based on the realization that there is no unique way to define an optimal aggregate for rankings, we investigate the predictive quality of a number of new aggregation methods and known fusion techniques from machine learning and social choice theory. b) Second, we propose a new approach to genomic data aggregation, termed *HyDRA* (Hybrid Distance-score Rank Aggregation), which combines the advantages of score-based and combinatorial aggregation techniques. We also propose incorporating a new *top-vs-bottom* (TvB) weighting feature into the hybrid schemes. The TvB feature ensures that aggregates are more reliable at the top of the list, rather than at the bottom, since only top candidates are tested experimentally. Specifically, we combine score-based Borda and Kendall permutation distance aggregation methods with TvB weightings. c) Third, we propose an iterative procedure for gene discovery that operates via successful augmentation of the set of training genes by genes discovered in previous rounds, checked for consistency.

We tested HyDRA on a number of gene sets, including Autism, Breast cancer, Colorectal cancer, Endometriosis, Ischaemic stroke, Leukemia, Lymphoma, and Osteoarthritis. Furthermore, we performed iterative gene discovery for Glioblastoma, Meningioma and Breast cancer, using a sequentially augmented list of training genes related to the Turcot syndrome, Li-Fraumeni condition and other diseases. The methods outperform state-of-the-art software tools such as ToppGene and Endeavour.

*To my parents.*

# ACKNOWLEDGMENTS

I would like to thank my dear adviser Professor Olgica Milenkovic for her support throughout my graduate studies. Particularly, I thank her for introducing me to the field of molecular biology and bioinformatics, providing the opportunity to apply tools from information theory and signal processing to solve problems in biology. Her careful guidance helped us to share the invaluable moments of an NSF fellowship award, HyDRA journal publication, and CSHL platform presentation. Without her, I would not have been who I am today. Likewise, I thank my collaborator Farzad Farnoud for insightful ideas and discussions, and my funding source the NSF Graduate Research Fellowship Program.

My second round of thanks goes to the mentors from UC San Diego who prepared me for graduate studies. I thank Professor Tara Javidi for a great introductory course on information and coding theory, offering an undergraduate research opportunity, and providing the last bits of encouragement for me to pursue graduate studies. I am grateful that Professor David Meyer introduced me to quantum computing and mentored me to complete my undergraduate honors thesis on quantum coding theory. I also thank Professor Paul Siegel for an excellent course on probabilistic coding theory and for providing me the opportunity to help him teach “ECE 101: linear systems” as a teaching assistant.

Lastly, I thank my father Hoyong Kim and my mother Namee Chung for moral and physical support through countless Kakaotalk messages, grandma-made Korean food packages, and many visits to Urbana-Champaign. I dedicate this thesis to my parents.

# CONTENTS

Chapter 1	INTRODUCTION . . . . .	1
1.1	Introduction to Molecular Biology . . . . .	1
1.2	Mutations and Diseases . . . . .	3
Chapter 2	MOTIVATION AND LITERATURE REVIEW . . . . .	7
2.1	Motivation . . . . .	7
2.2	Prior Work . . . . .	8
2.3	Our Contributions . . . . .	10
Chapter 3	SYSTEM AND METHODS . . . . .	12
3.1	Score-Based Methods . . . . .	13
3.2	Distance-Based Methods . . . . .	15
3.3	The Lovász-Bregman Divergence Method . . . . .	19
Chapter 4	ALGORITHMS AND IMPLEMENTATION . . . . .	22
4.1	Cross-Validation . . . . .	23
4.2	Gene Discovery . . . . .	23
Chapter 5	RESULTS . . . . .	26
5.1	Cross-Validation . . . . .	26
5.2	Gene Discovery . . . . .	30
Chapter 6	CONCLUSION . . . . .	34
	REFERENCES . . . . .	35

# Chapter 1

## INTRODUCTION

### 1.1 Introduction to Molecular Biology

We begin with a brief introduction to molecular biology, summarized from parts of Cooper and Hausman [1]. The basic cell theory states that 1) all living organisms are composed of one or more cells, 2) the cell is the most basic unit of life, and 3) all cells arise from pre-existing living cells. Thus, in order to study living organisms, such as humans, it is crucial to study our cells. The macromolecule deoxyribonucleic acid (DNA) contains the hereditary information that is passed on from cell to cell, thereby making it an important molecule to study. DNA is composed of two pyrimidines, cytosine (C) and thymine (T), and two purines, adenine (A) and guanine (G). The discovery of the structure of DNA goes back to 1949, when Chargaff realized that the amount of adenine was similar to the amount of thymine, and likewise the amount of guanine was almost the same as the amount of cytosine, which suggests A and T may be linked, and G and C may be related in structure as well. Later, Rosalind Franklin and Maurice Wilkins suggested a helical model of DNA, based on X-ray diffraction patterns, and James Watson and Francis Crick concluded that the only model that worked is a double helix (Figure 1.1). Today it is widely known that the DNA is a double helix with complementary base pairing, A with T and C with G.

DNA has the ability to replicate itself, thanks to the complementary base-pairing property. There were three proposed models for replication: semi-conservative, conservative, and dispersive, as shown in Figure 1.2. Our goal is to preserve the sequences of both strands of DNA after the replication, and the conservative model fails by preserving only one copy of itself and introducing an error prone replicate. Likewise, the dispersive model fails to produce even one copy of the original DNA. The most reasonable model

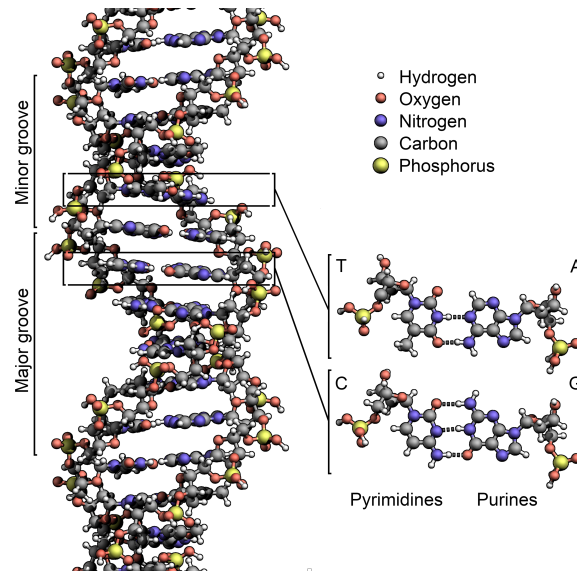


Figure 1.1: Complementary base-pairing of DNA. Source: wikipedia.

is semiconservative as it utilizes the complementary base pairing property. According to this model, one strand is used as a template to synthesize the complementary bases, and likewise the other strand is a template and adds on the complementary bases. As a result, the molecule successfully makes two copies of itself from one copy, where the red indicates one strand from the original copy and the pink corresponds to the newly synthesized strand. This process is known as *DNA Replication*.

Ribonucleic acid (RNA) is another information storage molecule made up of nucleotides, but unlike DNA, RNA does not have a hydroxyl group, and is usually single-stranded. RNA has two pyrimidines, thymine (T) and uracil (U), and two purines, adenine (A) and guanine (G). *Transcription* is the process of copying the DNA molecules into RNA molecules by using the complementary base-pairing property.

The information in RNA is used to build proteins, which performs many crucial functions in living organisms, through a process called *translation*. The protein is made up of a sequence of 20 amino acids according to the codon table in Table 1.1. As there are 4 bases in RNA, if two nucleotides code for one amino acid, there would be  $4^2 = 16$  amino acids, which is less than 20. Thus, each amino acid is made up of three nucleotides, and the genetic code is called a *triplet code*, and the nucleotides are read in groups of 3 called *codons*. This results in  $4^3 = 64$  possibilities, where 61 specify an



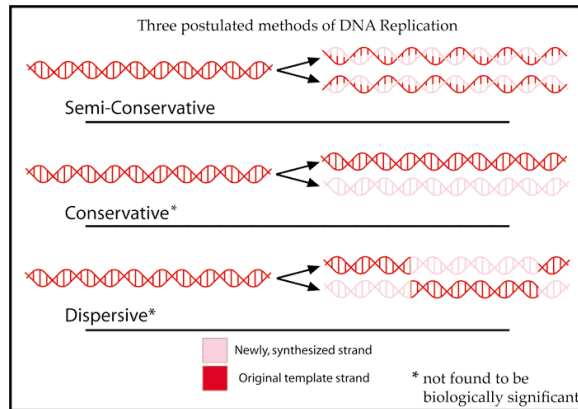


Figure 1.2: Semiconservative model of DNA replication. Source: wikipedia.

amino acid and 3 specify the end of the protein (STOP codon). As illustrated in Table 1.1, an amino acid can be specified by more than one codon, in which case the genetic code is considered *degenerate*.

As an analogy, DNA replication can be thought of as an oral statement of the sentence “Biology is fun”, the transcription is writing down the sentence in English, the only difference being the medium of language, and the translation is interpreting the sentence into a different language. The flow of genetic information from DNA to RNA to protein is known as the “Central Dogma of Molecular Biology”, as illustrated in Figure 1.3.

## 1.2 Mutations and Diseases

All three processes of the central dogma of molecular biology depend on reliable pairing of complementary bases. Unfortunately, no biological process is completely accurate in its implementation and execution. DNA replication is the most reliable of the three because mistakes in replication lead to alterations in the nucleotide sequence in the DNA and they are passed on to daughter cells when the cell divides. The heritable changes in the DNA are called *mutations*, and they can have positive or negative effects. Genes refer to DNA and RNA sequences that code for proteins and functional RNAs, which are responsible for genetic traits. As a general convention, we denote a normal, unmutated (wild type) gene by  $gene^+$  and a mutated form of the gene by  $gene^-$ . In single-cell organisms, all daughter cells have the mutation. In multi-cell organisms, mutations can be either somatic, that is passed to

	2nd base in the codon					
	U	C	A	G		
1st base in codon U	Phe	Ser	Tyr	Cys	U	3rd base in codon
	Phe	Ser	Tyr	Cys	C	
	Leu	Ser	STOP	STOP	A	
	Leu	Ser	STOP	Trp	G	
C	Leu	Pro	His	Arg	U	3rd base in codon
	Leu	Pro	His	Arg	C	
	Leu	Pro	Gln	Arg	A	
	Leu	Pro	Gln	Arg	G	
A	Ile	Thr	Asn	Ser	U	3rd base in codon
	Ile	Thr	Asn	Ser	C	
	Ile	Thr	Lys	Arg	A	
	Met	Thr	Lys	Arg	G	
G	Val	Ala	Asp	Gly	U	3rd base in codon
	Val	Ala	Asp	Gly	C	
	Val	Ala	Glu	Gly	A	
	Val	Ala	Glu	Gly	G	

Table 1.1: The genetic codon table.

daughter cells in the area, or germ-line, that is passed to new organisms. In general, germ-line mutations may affect the survivability of an organism, as all new cells acquire the mutation. On the other hand, mistakes in transcription or translation are not as critical because many copies of RNA are produced, and RNAs are not heritable over multiple generations.

Even though DNA corrects mistakes through proofreading and mismatch repair, it fails to correct all mistakes, leading to mutations carried over to daughter cells. There are 4 main types of mutations: base substitutions, frameshift, insertion, and deletion. Base substitutions and frameshift mutations are referred to as point mutations, and insertions and deletions (which are sometimes called “indels” collectively) are on a larger scale, chromosomal mutations.

Base substitution mutations occur when the DNA polymerase mistakenly replaces one basepair with another. There are 3 categories of base substitution mutations: missense, nonsense, same sense. Suppose the template DNA was AAC, which codes for the UUG = Leu protein. If the first base changed to a C, which codes for GUG = Val, a missense mutation occurred, where the message has been changed. If the second base changed to a T,

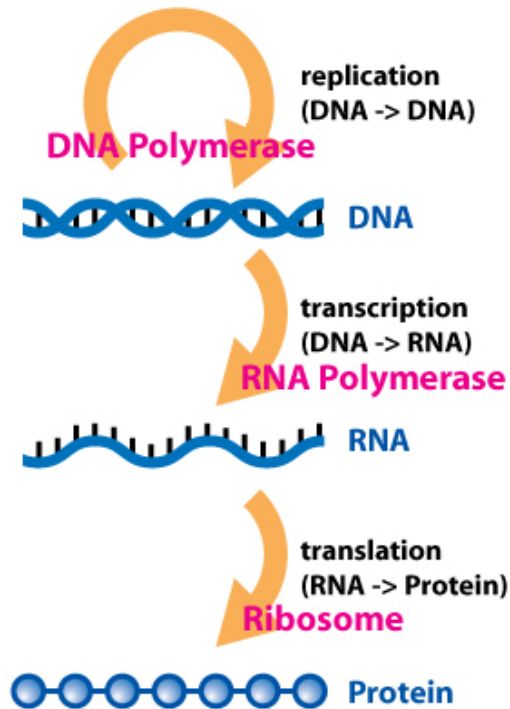


Figure 1.3: Central dogma of molecular biology. Source: wikipedia.

which codes for UAG = STOP, this truncates the message due to the stop codon. On the other hand, if the third base changed to a T, encoding UUA = Leu, the message did not change at all, and the same sense mutation occurred. Same sense mutations may not have phenotypical consequences, as they do not change the encoded protein. Nevertheless, the full impacts of sense mutations are currently not well understood. Nonsense mutation causes premature termination and may affect the protein function if the region was important. However, missense mutation can potentially be disastrous. For example, sickle cell anemia, which is characterized by defective  $\beta$ -globin subunit in hemoglobin protein, is caused by a change in the  $\beta$ -globin gene at position 7. Likewise, phenylketonuria (PKU) is also caused by a single base change in an enzyme.

Frameshift mutations disrupt the normal reading frame by an insertion or deletion of 1 or 2 bases. As a result, all codons after the insertion or deletion are altered and usually result in premature truncation. On the other hand, if there is a deletion of 3 bases, the protein only loses one amino acid and does not change the reading frame, making it non-frameshift mutation. However, the loss of a single amino acid could be detrimental as well. Cystic fibrosis

is characterized by faulty protein *CFTR*, where 70% of the cystic fibrosis is caused by an in-frame removal of a Phe codon in *CFTR*, preventing the protein from folding properly.

At the chromosomal level, there are two main types of mutations: insertions and deletions. Insertions are caused by large pieces of DNA inserted into gene sequence, usually more than 100 bases. Examples of diseases caused by insertions are: 1) Huntington's disease, where the *HTT* gene contains a repeated sequence of CAG; 2) hemophilia, caused by insertions in the *F8* gene; and 3) fragile X syndrome, a result of CGG repeats in the *fragile X mental retardation 1 (FMR1)* gene on the X chromosome. On the other hand, deletions are due to loss of large pieces of DNA sequence. For instance, Duchenne's muscular dystrophy disorder is caused by deletions in the *dystrophin (DMD)* gene.

All of the above relationships between genetic mutations and diseases have been studied through long, laborious "knockout" experiments, where the scientists genetically modify a mouse and observe changes in its behavior and appearance. As data scientists, we propose a solution, collectively called "gene prioritization," to eliminate the number of knockout experiments, thereby saving time, cost, and mice. In the subsequent chapters, we delve into gene prioritization methods, which have been previously presented in [2].

## Chapter 2

# MOTIVATION AND LITERATURE REVIEW

### 2.1 Motivation

Fundamental results from social choice theory, political and computer sciences, and statistics have shown that there exists no consistent, fair and unique way to aggregate rankings. Instead, one has to decide on an aggregation approach using a predefined set of desirable properties for the aggregate. The aggregation methods fall into two categories, score-based and distance-based approaches, each of which has its own drawbacks and advantages. This work is motivated by the observation that merging these two techniques in a computationally efficient manner, and by incorporating additional constraints, one can ensure that the predictive quality of the resulting aggregation algorithm is very high. We refer to these methods as hybrid methods, and outline how they can be used in an iterative discovery procedure and subsequently aggregated to produce very accurate gene rankings.

Identification of genes that predispose an individual to a disease is a problem of great interest in medical sciences and systems biology [3]. The most accurate and powerful methods used for identification are experimental in nature, involving normal and disease samples [4]. Experiments are time-consuming and costly, complicated by the fact that typically, multiple genes have to be jointly mutated to trigger the onset of a disease. Given the large number of human genes ( $\geq 25,000$ ), testing even relatively small subsets of pairs of candidate genes is prohibitively expensive [5].

To mitigate this issue, a set of predictive analytical and computational methods have been proposed under the collective name *gene prioritization techniques*. Gene prioritization refers to the complex procedure of ranking genes according to their likelihoods of being linked to a certain disease. The likelihood function is computed based on multiple sources of evidence, such

as sequence similarity, linkage analysis, gene annotation, functionality and expression activity, gene product attributes – all determined with respect to a set of training genes.

## 2.2 Prior Work

A wide range of tools has been developed for identifying genes involved in a disease, as surveyed [6]. Existing software includes techniques based on network information, such as GUILDify [7] and GeneMANIA [8], data mining and machine learning-based approaches as described in [9], POCUS [10], and SUSPECTS [3], and methods using statistical analysis, including Endeavour [11, 12], ToppGene [13], and NetworkPrioritizer [14]. Here, we focus on statistical approaches coupled with new combinatorial algorithms for gene prioritization, and emphasize one aspect of the prioritization procedure: rank aggregation.

The methods of choice for aggregating multiple sources of similarity evidence include  $Q$ -statistics and order statistics (Endeavour) and Fisher statistics (ToppGene). An exception is the very recent work NetworkPrioritizer, which relies on a small number of well known *combinatorial techniques* from social choice theory that does not make direct use of scores in the form of p-values.

The problem of aggregating rankings of distinct objects or entities provided by a number of experts, voters, or search engines has a rich history [15]. One of the key findings is that various voting paradoxes arise when more than three candidates are to be ranked: it is frequently possible not to have a candidate that wins all pairwise competitions (the Condorcet paradox) and it is theoretically impossible to guarantee the existence of an aggregate solution that meets a certain predefined set of criteria (such as those imposed by Arrow’s impossibility theorem [15]). These issues carry over to aggregation methods used for gene discovery, and as a result, the rank-ordered lists of genes heavily depend on the particular aggregation method used.

Two families of methods have found wide applications in rank aggregation: combinatorial methods (including score- and distance-based approaches) [16] and statistical methods. In the bioinformatics literature, the aggregation methods of choice are statistical in nature, relying on pre-specified hypothe-

ses to evaluate the distribution of the gene rankings. One of the earliest prioritization softwares, Endeavour, uses the  $Q$ -statistics for multiple significance testing, and measures the minimum false discovery rate at which a test may be called significant. In particular, rankings based on different similarity criteria are combined via order statistics approaches. For this purpose, one uses the rank ratio (normalized ranking) of a gene  $g$  for  $m$  different criteria,  $r_1(g), \dots, r_m(g)$  and recursively computes the  $Q$ -value, defined as

$$Q_i(r_1(g), \dots, r_m(g)) = m! \int_0^{r_1(g)} \int_{s_1}^{r_2(g)} \dots \int_{s_{m-1}}^{r_m(g)} d_{s_m} d_{s_{m-1}} \dots d_{s_1}.$$

Post-processed  $Q$ -values are used to create the resulting ranking of genes. The drawbacks of the method are that it is based on a null hypothesis that is difficult to verify in practice, and that it is computationally expensive, as it involves evaluating an  $m$ -fold integral. To enable efficient scaling of the method, Endeavour resorts to approximating the  $Q$ -integral. The influence of the approximation errors on the final ranking is hard to assess, as small changes in scores may result in significant changes of the aggregate orderings.

Likewise, ToppGene uses a well-known statistical approach, called the Fisher  $\chi^2$  method. It first determines the p-values of similarity score indexed by  $j$ , denoted by  $p(j)$ , for  $j = 1, \dots, m$ . The p-values are computed through multiple pre-processing stages, involving estimation of the information contents (i.e., weights) of annotation terms, setting up a similarity criteria based on Sugeno fuzzy measures (i.e., non-additive measures) [17], and performing meta-testing. The use of fuzzy measures ensures that all similarities are non-negative. Then, under the hypothesis of *independent tests*, ToppGene uses Fisher's inverse  $\chi^2$  result, stating that  $-2 \sum_{j=1}^m \log p(j) \rightarrow \chi^2(2m)$ . Here,  $\chi^2(2m)$  stands for the chi-square distribution with  $2m$  degrees of freedom. The result is asymptotic in nature, and based on possibly impractical and unverifiable independence assumptions.

A number of methods, and additive scoring methods in particular, have the following drawbacks: a) they assume that only the total score matters, and the balance between the number of criteria that highly ranked the gene and those that ranked it very low is irrelevant. For example, outlier rankings may reduce the overall ranking of a gene to the point that it is not considered a disease gene candidate, while the outlier itself may be a problematic criterion.

To illustrate this observation, consider a gene that was ranked 1st, 2nd, 1st, 20th by four criteria. At the same time, consider another gene that was ranked 6th by all four criteria. It may be unclear which of these two genes is more likely to be involved in the disease, given that additive score methods would rank the two genes equally (as one has  $(1+2+1+20)/4=6$ ). Nevertheless, it appears reasonable to assume that the first candidate is a more reliable choice for a disease gene, as it had a very high ranking for three out of four criteria. b) No distinction is made about the accuracy of ranking genes in any part of the list; i.e., the aggregate ranking has to be *uniformly accurate* at the top, middle and bottom of the list. Clearly, neither of the two aforementioned assumptions is justified in the gene prioritization process: there are many instances where genes similar only under a few criteria (such as sequence similarity or linkage distance) are involved in the same disease pathway. Furthermore, as the goal of prioritization is to produce a list of genes to be experimentally tested, only the highest ranked candidate genes are important and should have higher accuracy than other genes in the list. In addition, most known aggregation methods are highly sensitive to outliers and ranking errors.

## 2.3 Our Contributions

We propose a new approach to gene prioritization by introducing a number of novel aggregation paradigms, which we collectively refer to as *HyDRA* (Hybrid Distance-score Rank Aggregation). The gist of HyDRA is to combine *combinatorial approaches* that have universal axiomatic underpinnings with *statistical evidence* pertaining to the accuracy of individual rankings. Our preferred distance measure for combinatorial aggregation is the Kendall distance [18], which counts the number of pairwise disagreements between two rankings, and was axiomatically postulated by Kemeny in [16]. The Kendall distance is closely related to the Kendall rank correlation coefficient [19], [20]. As such, it has many properties useful for gene prioritization, such as monotonicity, reinforcement and Pareto efficiency [21]. The Kendall distance can be generalized to take into account positional relevance of items, as was done in our companion paper [22]. There, it was shown that by assigning weights to pairs of positions in rankings, it is possible to a) eliminate negative outliers



from the aggregation process, b) include quantitative data into the aggregate, and c) ensure higher accuracy at the top of the ranking than at the bottom.

The contributions of this work are threefold. First, we introduce new weighted distance measures, where we compute the weights based on statistical evidence in the form of a function of the difference between p-values of adjacently ranked items. Aggregation weights based on statistical evidence improve the accuracy of the combinatorial aggregation procedure and make the aggregate more robust to estimation errors. Second, we describe how to scale the weights obtained based on statistical evidence by a decreasing sequence of TvB (Top versus Bottom) multipliers that ensure even higher accuracy at the top of the aggregated list. As aggregation under the Kendall metric is NP hard [23], and the same is true of the weighted Kendall metric, we propose a 2-approximation method that is stable under small perturbations. Aggregation is accomplished via weighted bipartite matching, such as the Hungarian algorithm and derivatives thereof [24]. Third, we test HyDRA within two operational scenarios: cross-validation and disease gene discovery. In the former case, we assess the performance of different hybrid methods with respect to the choice of the weighting function and different number of test and training genes. In the latter case, we adapt aggregation methods to gene discovery via a new iterative re-ranking procedure.

## Chapter 3

# SYSTEM AND METHODS

In our subsequent exposition, we use Greek lower case letters to denote complete linear orders (permutations), and unless explicitly mentioned otherwise, our findings also hold for partial (incomplete) permutations. Latin lower case letters are reserved for score vectors or scalar scores, and which of these entities we refer to will be clear from the context. The number of test genes equals  $n$ , while the number of similarity criteria equals  $m$ . Throughout the chapter, we also use  $[k]$  to denote the set  $\{1, \dots, k\}$  and  $\mathbb{S}_n$  to denote the set of all permutations on  $n$  elements – the symmetric group of order  $n!$ .

For a permutation  $\sigma = (\sigma(1), \dots, \sigma(n))$ , the rank of element  $i$  in  $\sigma$ ,  $\text{rank}_\sigma(i)$ , equals  $\sigma^{-1}(i)$ , where  $\sigma^{-1}$  denotes the inverse permutation of  $\sigma$ . For a vector of scores  $x = (x(i))_{i=1}^n \in \mathbb{R}^n$ ,  $\sigma_x$  represents a permutation describing the scores in decreasing order, i.e.,  $\sigma_x(i) = \text{argmax}_{k \in T_i} x(k)$ , where  $T_i$  is defined recursively as  $T_i = T_{i-1} \setminus \sigma_x(i)$ , with  $T_0 = [n]$ . For example, if  $x = (2.5, 3.8, 1.1, 0.7)$ , then  $\sigma_x = (2, 1, 3, 4)$ . Note that if  $p$  is a vector of p-values, higher scores are associated with smaller p-values, so that  $\text{argmax}$  should be replaced by  $\text{argmin}$ .

The terms *gene* and *element* are used interchangeably, and each permutation is tacitly assumed to be produced by one similarity criterion. For a set of permutations  $\Sigma = \{\sigma_1, \dots, \sigma_m\}$ ,  $\sigma_i = (\sigma_i(1), \dots, \sigma_i(n))$ , an *aggregate permutation*  $\sigma^*$  is a permutation that optimally represents the rankings in  $\Sigma$ . Combinatorial aggregates may be obtained using *score-based* and *distance-based* methods. Note that score and distance-based methods do not make use of quantitative information such as, for example, p-values (for the case of gene prioritization) or ratings (for the case of social choice theory and recommender systems). In what follows, we briefly describe score and distance-based methods and introduce their *hybrid* counterparts which allow to integrate p-values and relevance constraints into combinatorial aggregation approaches.

### 3.1 Score-Based Methods

Score-based methods are the simplest and computationally least demanding techniques for rank aggregation. As inputs, they take a set of permutations or partial permutations,  $\Sigma = \{\sigma_1, \dots, \sigma_m\}$ ,  $\sigma_i = (\sigma_i(1), \dots, \sigma_i(n))$ . For each permutation  $\sigma_i \in \Sigma$ , the scoring rule awards  $s(\sigma_i(1), i)$  points to element  $\sigma_i(1)$ ,  $s(\sigma_i(2), i)$  points to element  $\sigma_i(2)$ , and so on. For a fixed  $i$ , the scores are non-increasing functions of their first index. Each element  $k \in [n]$  is assigned a cumulative score equal to  $\sum_{j=1}^m s(k, j)$ . The simplest scoring method is Borda's count, for which  $s(k, j) = n - k + 1$  independent on  $j$ .

The Borda count and related scoring rules exclusively use positional information in order to provide an aggregate ranking. Ignoring actual p-values (ratings) may lead to aggregation problems, as illustrated by the next example.

**Example 1:** Assume that  $n = 5$  elements were rated according to  $x = (7.0, 7.01, 0.2, 0.45, 7.001)$ . The ranking induced by this rating equals  $\sigma_x = (2, 5, 1, 4, 3)$ , indicating that element 2 received the highest rating, element 5 received the second highest rating and so on. According to the Borda rule, element 2 receives 5 points, element 5 receives 4 points, etc. Despite the fact that candidates 2 and 1 are almost tied with scores of 7.01 and 7.0, and that the difference in their scores may be attributed to computational imprecision, element 2 receives 5 points while element 1 receives only 3 points. As a result, very small differences in ratings may result in large differences in Borda scores.

One way to approach the problem is to quantize the score and work with rankings with ties, instead of full linear orders (i.e. permutations). Elements tied in their rank receive the same number of points in the generalized Borda scheme. A preferred alternative, which we introduce in this work, is the *Hybrid Borda method*.

Let  $p(i, j)$  denote the p-value of gene  $i$  computed under similarity criteria  $j, j = 1, \dots, m$ . The cumulative score of element  $i$  in the hybrid Borda setting

is computed as

$$S_i = \sum_{j=1}^m \left( \frac{\sum_{k \neq i} p(k, j) \mathbb{1}_{\{p(k, j) \geq p(i, j)\}}}{p(i, j)} \right).$$

The overall aggregate is obtained by ordering  $S$  in a descending order. It is straightforward to see that the previous score function extends the Borda method in so far as it scores an element (gene) according to the total score of elements ranked lower than the element. Recall that in Borda's method, the element ranked  $i$  is awarded  $n - i + 1$  points, as  $n - i + 1$  elements are ranked below it, each receiving the same score 1. In our Hybrid Borda method, each element is awarded a score in *accordance with the  $p$ -values of elements ranked below it*.

**Example 2:** Let  $n = 4$  and  $m = 2$ , where the two ratings equal to  $p_1 = (0.2, 0.3, 0.01, 0.12)$  and  $p_2 = (0.1, 0.4, 0.2, 0.35)$ . The Hybrid Borda scores  $S_i$  for genes  $i = 1, 2, 3, 4$  are computed as  $S_1 = \frac{0.3}{0.2} + \frac{(0.4+0.2+0.35)}{0.1} = 11$ ,  $S_2 = 0$ ,  $S_3 = \frac{(0.2+0.3+0.12)}{0.01} + \frac{(0.4+0.35)}{0.2} = 65.75$ , and  $S_4 = \frac{(0.2+0.3)}{0.12} + \frac{0.4}{0.35} = 5.3$ . By ordering the values  $S_i$  in a descending manner, we obtain the overall aggregate  $\sigma_{HB} = (3, 1, 4, 2)$ .

The hybrid Borda method can be extended further by adding a TvB feature, resulting in the *Weighted* Hybrid Borda method. This is accomplished by including *increasing* (multiplier) weights into the score aggregates, thus stressing the top of the list more than the bottom. More precisely, the score of gene  $i$  is computed as:

$$S_i = \sum_{j=1}^m \left( \frac{\sum_{k \neq i} wm(k, j) p(k, j) \mathbb{1}_{\{p(k, j) \geq p(i, j)\}}}{wm(i, j) p(i, j)} \right),$$

where one simple choice for the weight multipliers that provides good empirical performance equals

$$wm(i, j) = \frac{1}{n - \text{rank}_{\sigma_j}(i) + 1}.$$

Note that other weight functions are possible as well, but we used the above formula for its simplicity and good empirical performance.

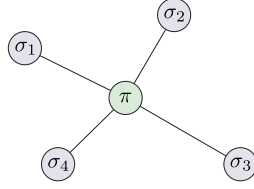


Figure 3.1: Four rankings:  $\sigma_1, \sigma_2, \sigma_3, \sigma_4$ , and their aggregate (median) ranking  $\pi$ .

## 3.2 Distance-Based Methods

Another common approach to rank aggregation is distance-based rank aggregation. As before, assume that one is given a set of permutations  $\Sigma = \{\sigma_1, \dots, \sigma_m\}$ . For a given distance function between two permutations  $\sigma$  and  $\pi$ ,  $d(\sigma, \pi)$ , aggregation reduces to

$$\pi = \arg \min_{\sigma} \sum_{i=1}^m d(\sigma, \sigma_i).$$

The aggregate  $\pi$  is frequently referred to as the *median* of the permutations, and is illustrated in Figure 3.1.

Table 3.1: Two frequently used distance measures for permutations, accounting for swaps or element-wise differences. In the second example, the Kendall tau distance between the permutation  $\sigma_1 = (a, b, c)$  and  $\sigma_2 = (c, b, a)$  equals 3: one first swaps elements at positions 1 and 2 to get  $(b, a, c)$ , then elements at positions 2 and 3 to get  $(b, c, a)$ , and finally elements at positions 1 and 2 to get  $\sigma_2 = (c, b, a)$ . All swaps contribute the same weight (one) to the distance.

Distance	Measurement	Example
Spearman's footrule	Sum of differences of ranks of elements.	$d_F(abc, cba) = 2 + 0 + 2 = 4$
Kendall	Minimum number of adjacent swaps of entries for transforming one ranking into another.	$d_K(abc, cba) = 3$

One of the most important features of distance-based approaches is the choice of the distance function. Table 3.1 lists two of the most frequently used distances, the Kendall tau distance and the Spearman footrule. As may be seen from the table, the distance measures are combinatorial in nature,

and do not account for scores or p-values. Furthermore, as already mentioned in the introduction, it is known that aggregation under the Kendall metric is computationally hard. Nevertheless, there exist a number of techniques which provide provable *approximation* guarantees for the aggregate, including the weighted Bipartite Graph Matching (WBGGM) method (using the fact that the Spearman distance aggregate is a 2-approximation for the Kendall aggregate), linear programming (LP) relaxation, and Page Rank/Markov chain (PR) methods [20, 25, 26].

The Kendall distance also does not take into account the fact that the top of a list is more important than the remainder of the list. To overcome this problem, we introduced the notion of *weighted Kendall distances*, where each adjacent swap is assigned a cost, and where the cost is higher at the top of a list. This ensures that in an aggregate, strong showings of candidates are emphasized compared to their weaker showings, accounting for the fact that it is often sufficient to have strong similarity with respect to only a subset of criteria. Furthermore, such weights ensure that higher importance is paid to the top of the aggregate ranking.

The idea behind the weighted Kendall distance  $d_w$  is to compute this distance as the shortest path in a graph describing swap relationships between permutations. The key concepts are illustrated in Figures 3.2 and 3.3, where each edge is assigned a length proportional to its weight  $W$ . This weight depends on the swap being made at the top or at some other position in the ranking. Given that it is computationally demanding to aggregate under the weighted Kendall distance, we use a specialized approximation function  $D_w(\sigma, \theta)$  for  $d_w$ , of the form

$$D_w(\sigma, \theta) = \sum_{i=1}^n w(\sigma^{-1}(i) : \theta^{-1}(i)), \quad (3.1)$$

where

$$w(k : l) = \begin{cases} \sum_{h=k}^{l-1} W(h, h+1), & \text{if } k < l, \\ \sum_{h=l}^{k-1} W(h, h+1), & \text{if } k > l, \\ 0, & \text{if } k = l \end{cases} \quad (3.2)$$

denotes the sum of the weights of edges  $W(\cdot)$  representing adjacent transpositions  $(k \ k+1), (k+1 \ k+2), \dots, (l-1 \ l)$ , if  $k < l$ , the sum of the weights of edges  $W(\cdot)$  representing adjacent transpositions  $(l \ l+1), (l+1 \ l+2), \dots, (k-1 \ k)$ ,

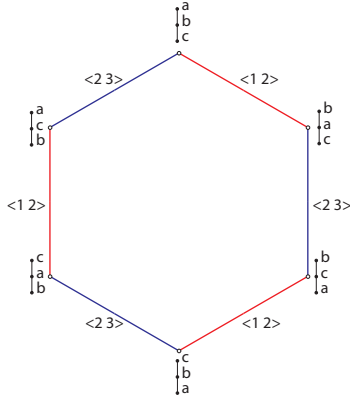


Figure 3.2: The Kendall distance is the weight of the shortest path between two vertices labeled by two permutations, with each edge having length (weight) one. Edges are labeled by the adjacent swaps used to move between the vertex labels. For example, the two vertices labeled by  $acb$  and  $cab$  are connected via an edge bearing the label  $\langle 12 \rangle$ , indicating that the two permutations differ in one swap involving the first and second element.

if  $l < k$ , and 0, if  $k = l$ .

**Example 3:** Suppose that one is given four rankings,  $(1, 2, 3)$ ,  $(1, 2, 3)$ ,  $(3, 2, 1)$  and  $(2, 1, 3)$ . There are two optimal aggregates according to the Kendall  $\tau$  distance, namely  $(1, 2, 3)$  and  $(2, 1, 3)$ . Both have cumulative distance four from the set of given permutations. If the transposition weights are non-uniform, say such that  $W(12) > W(23)$ , the solution becomes unique and equal to  $(1, 2, 3)$ . If the last ranking is changed from  $(2, 1, 3)$  to  $(2, 3, 1)$ , exactly three permutations are optimal from the perspective of Kendall  $\tau$  aggregation:  $(1, 2, 3)$ ,  $(2, 1, 3)$ , and  $(2, 3, 1)$ . These three solutions give widely different predictions of what one should consider the top candidate. Nevertheless, by choosing once more  $W(12) > W(23)$  the solution becomes unique and equal to  $(1, 2, 3)$ .

It can be shown that for any non-negative weight function  $w$ , and for two permutations  $\sigma$  and  $\theta$ , one has

$$1/2D_w(\sigma, \theta) \leq d_w(\pi, \sigma) \leq D_w(\sigma, \theta).$$

In a companion paper [25], we presented extensions of the WBGm and PR aggregation methods for weighted Kendall distances. Here, we will pursue the WBGm framework, and propose a new method to compute the weights  $W(\cdot)$  of edges (swaps) based on the p-values of the genes within each similarity

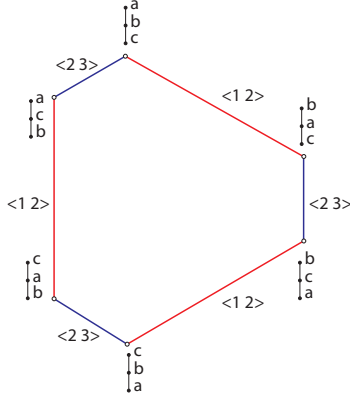


Figure 3.3: The weighted Kendall distance is the weight of the shortest path between two permutations, with edges having possibly different lengths (weights). Edges are labeled by the adjacent swaps used to move along the vertices.

criterion ranking. We refer to the resulting weighted model as the *Hybrid Kendall* method.

To start, arrange the p-values of all genes based on all similarity criteria into an  $n \times m$  matrix  $P$ . Next, rearrange the p-values of genes for each criterion in an increasing order, and denote the resulting rearranged matrix by  $P^* = (p^*(i, j))$ . We use the following  $(n - 1) \times m$  swap weight matrix  $\mathcal{W}$ , with entries

$$\mathcal{W}(i, j) = c \left( \frac{P^*(i + 1, j) - P^*(i, j)}{P^*(i + 1, j)} \right) \times d^{n-i},$$

indicating how much it costs to swap positions  $i$  and  $i + 1$  for criterion  $j$ . The parameters  $c, d$  are constants independent of  $n$  and  $m$ , used for normalization and for emphasizing the TvB constraint, respectively. For our simulations, we set  $c = 10$  and  $d = 1.05$ , as these choices provided good empirical performance on synthetic data. The swap matrix assigns high weight to the top of the list.

To compute the aggregate based on the approximate distance  $D_w(\theta, \sigma)$ , we only need to accumulate each of the contributions from the training permutations in  $\Sigma$ . This may be achieved by using a  $n \times n$  total cost matrix  $C$ , with entry  $C(i, j)$  indicating how much it would “cost” for gene  $i$  to be



ranked at position  $j$ :

$$C(i, j) = \frac{1}{m} \sum_{k=1}^m \sum_{l=\min(j, \sigma_{p_k}(i))}^{\max(j, \sigma_{p_k}(i))-1} \mathcal{W}(l, k).$$

The total cost matrix  $C$  is the input to the WBGGM algorithm, where  $C(i, j)$  denotes the weight of an edge connecting gene  $i$  with position  $j$  (see Figure 3.4 for an example of the bipartite graph, with the left-hand side nodes denoting genes and the right-hand side nodes denoting their possible positions; the minimum weight matching is represented by bold edges). To find the minimum cost solution, or the maximum weight matching, we used the classical *Hungarian algorithm* [24] implemented in [27].

**Example 4:** Let  $n = 4$  and  $m = 2$ , where the two ratings equal to  $p_1 = (0.2, 0.3, 0.01, 0.12)$  and  $p_2 = (0.1, 0.4, 0.2, 0.35)$ . Then

$$P^* = \begin{bmatrix} 0.01 & 0.1 \\ 0.12 & 0.2 \\ 0.2 & 0.35 \\ 0.3 & 0.4 \end{bmatrix}, \quad \mathcal{W} = \begin{bmatrix} 10.61 & 5.79 \\ 4.41 & 4.73 \\ 3.5 & 1.31 \end{bmatrix},$$

$$C = \begin{bmatrix} 7.51 & 5.1 & 5.23 & 7.67 \\ 15.18 & 6.98 & 2.4 & 0 \\ 2.9 & 5.3 & 9.88 & 12.28 \\ 10.57 & 2.37 & 2.2 & 4.61 \end{bmatrix}.$$

For example, since gene 3 was ranked 1st and 2nd by the two criteria,  $C(3, 3) = \frac{1}{2}(10.61 + 4.41) + \frac{1}{2}(4.73) = 9.88$ . The minimum cost solution of the matching with cost matrix  $C$ , based on the Hungarian algorithm, yields the aggregate  $\sigma_{HK} = (3, 1, 4, 2)$ .

### 3.3 The Lovász-Bregman Divergence Method

A previously reported, distance measure represents another possible mean for performing hybrid rank aggregation. The so called *Lovász-Bregman* method [28] calls for a distance measure between real-valued vectors  $x \in \mathbb{R}_{\geq 0}^n$  and permutations.

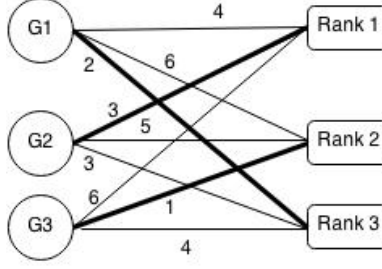


Figure 3.4: A matching in a weighted bipartite graph.

To define the Lovász-Bregman divergence that acts as a distance proxy between rankings and ratings, we start with a submodular set-function, i.e. a function  $f$  such that for a finite ground set  $V$ ,  $f : 2^V \rightarrow \mathbb{R}$ , and for all  $S, T \subset V$ , it holds  $f(S) + f(T) \geq f(S \cup T) + f(S \cap T)$ . The Lovász extension of  $f$ ,  $f^L(x)$ , equals

$$f^L(x) = \sum_{i=1}^n x(\sigma_x(i)) [f(S_j^{\sigma_x}) - f(S_{j-1}^{\sigma_x})],$$

where  $S_j^{\sigma_x}$  denotes the set  $\{\sigma_x(1), \dots, \sigma_x(j)\}$ . Note that under some mild conditions, the Lovász extension is convex. Let us next define the differential of  $f$  as

$$h_{\sigma_x}^f(\sigma_x(j)) = f(S_j^{\sigma_x}) - f(S_{j-1}^{\sigma_x}).$$

Then the Lovász-Bregman divergence is defined via the dot product

$$d_r(x||\sigma) = x \cdot (h_{\sigma_x}^f - h_{\sigma}^f).$$

Despite its seemingly complex expression, the Lovász-Bregman divergence allows for *closed form* aggregation for a large class of submodular functions  $f$ . The optimal aggregate reduces to the ranking induced by the *sum of real-valued rating vectors*, ordered in a decreasing manner.

If, as before,  $p(i, j)$  denotes the p-value of gene  $i$  under criterion  $j$ , we define the *normalized* Lovász-Bregman score for gene  $i$  as

$$\mathcal{L}(i) = \sum_{j=1}^m \frac{p(i, j)}{\frac{1}{n} \sum_{i=1}^n p(i, j)},$$

where the sum of p-values over criteria is normalized by the average of the p-values for each criterion. The aggregate equals  $\sigma_{\mathcal{L}}$ , where  $\mathcal{L} = (\mathcal{L}(i))_{i=1}^n$ .

**Example 5:** Let  $n = 4$  and  $m = 2$ , where the two ratings equal to  $p_1 = (0.2, 0.3, 0.01, 0.12)$  and  $p_2 = (0.1, 0.4, 0.2, 0.35)$ . Note that  $\frac{1}{n} \sum_{i=1}^n p(i, 1) = \frac{1}{4}(0.2 + 0.3 + 0.01 + 0.12) = 0.1575$ , and  $\frac{1}{n} \sum_{i=1}^n p(i, 2) = \frac{1}{4}(0.1 + 0.4 + 0.2 + 0.35) = 0.2625$ . The Lovász-Bregman scores  $\mathcal{L}(i)$ ,  $i = 1, 2, 3, 4$ , equal  $\mathcal{L}(1) = \frac{0.2}{0.1575} + \frac{0.1}{0.2625} = 1.65$ ,  $\mathcal{L}(2) = \frac{0.3}{0.1575} + \frac{0.4}{0.2625} = 3.43$ ,  $\mathcal{L}(3) = \frac{0.01}{0.1575} + \frac{0.2}{0.2625} = 0.83$ ,  $\mathcal{L}(4) = \frac{0.12}{0.1575} + \frac{0.35}{0.2625} = 2.1$ . By ordering  $\mathcal{L}(i)$  in an ascending manner, one arrives at  $\sigma_{LB} = (3, 1, 4, 2)$ .

## Chapter 4

# ALGORITHMS AND IMPLEMENTATION

We now turn our attention to testing different aggregation methods on lists of p-values generated by Endeavour and ToppGene. The aforementioned methods rely on a set of training genes known to be involved in a disease. The test genes are compared to all the training genes according to a set of similarity criteria, and the p-value of each comparison is computed in the process. For example, if the criterion is sequence similarity, the p-value reflects the z-value, describing the number of standard deviations above the mean for a given observation. Given the p-values, the question of interest becomes how to aggregate them into one ranking. Computing the p-values is a routine procedure, and the challenge of the prioritization process is to most meaningfully and efficiently perform the aggregation step.

There are two settings in which one can use the aggregation algorithms. The first setting is *cross-validation*, a verification step that compares the output of an aggregation algorithm with existing, validated knowledge. This mode of operation is aimed at discovering shortcomings and advantages of different methods. In the second setting, termed *gene discovery*, the aim is to identify sets of genes implicated in a disease which are not included in the database. Clearly, cross-validation studies are necessary first steps in gene discovery procedures, as they explain best aggregation strategies for different datasets and different similarity and training conditions.

For both methods, a list of genes involved in a certain disease (referred to as onset genes) was obtained from the publicly available databases OMIM [29] and/or the Genetic Association Database (GAD) [30]. Both of these sources rely on the literature about genetic association for a vast number of diseases, but OMIM typically provides a more conservative (i.e., shorter) list than the GAD. Onset genes were tested along with *random genes*, obtained by randomly permuting 19,231 human genes in the GeneCards database [31], and retaining the top portion of the list according to the chosen number of

test genes.

## 4.1 Cross-Validation

We performed a systematic, comparative performance analysis of the ToppGene and Endeavour aggregation algorithms and the newly proposed hybrid methods. Given a list of  $r$  onset genes, we first selected  $t$  onset genes to serve as target genes (henceforth referred to as *target onset genes*) for validation; we used the remaining  $r - t$  onset genes as training genes. Of the  $n$  test genes,  $n - t$  genes were selected randomly from GeneCards [31]. Our cross-validation procedure closely followed that of Endeavour and ToppGene: we fixed  $t = 1$ , and tested *all*  $r$  individual genes from the pool of onset genes, and then averaged the results. Averaging was performed as follows: we took target onset genes one-by-one and averaged their rankings over  $\binom{r}{t}_{t=1} = r$  experiments. Note that in principle, one may also choose  $t \geq 2$ ; in this case, the lowest ranking of the  $t$  genes (i.e., the highest positional value that a target onset gene assumed) should serve as a good measure of performance. One would then proceed to average the resulting rankings over  $\binom{r}{t}$  experiments, producing a “worst case scenario” for ranking of target onset genes. For fair comparison with Endeavour and ToppGene, we only used the first described method with  $t = 1$  and the same set of p-values as inputs. As will be described in subsequent sections, we used  $t \geq 2$  for gene discovery procedures.

## 4.2 Gene Discovery

The ultimate goal of gene prioritization is to *discover* genes that are likely to be involved in a disease without having any prior experimental knowledge about their role. We describe next a new, iterative *gene discovery* method. The method uses aggregation techniques or combinations of aggregation techniques deemed to be most effective in the cross-validation study.

Given a certain disease with  $r$  onset genes, we first identify  $s$  *suspect genes*. Suspect genes are genes that are known to be involved in diseases *related to*

that under study,<sup>1</sup> but have not been tested in this possible role. Suspect genes are processed in an iterative manner, as illustrated in Algorithm 1. In the first iteration,  $r$  onset genes are used for training, and  $s$  suspect genes, along with  $n - s$  randomly selected genes, are used as test genes. From the aggregate results provided by different hybrid algorithms, we selected  $q$  top-ranked genes and moved them to the set of training genes and simultaneously declared them as potential disease genes. The choice for the parameter  $q$  is governed by the number of training and test genes, as well as the empirical performance of the aggregation methods observed during multiple rounds of testing. The second iteration starts with  $r + q$  training genes,  $s - q$  suspect genes, and  $n - s + q$  randomly selected genes; the procedure is repeated until a predetermined stopping criterion is met, such as the size of the set of potential disease genes exceeding a given threshold.

---

<sup>1</sup>As an example, a suspect gene for glioblastoma may be a gene known to be implicated in another form of brain cancer, say meningioma.

---

**Algorithm 1:** Gene Discovery

---

**Input:** Set of onset genes,  $O = \{o_1, o_2, \dots, o_r\}$ , set of suspect genes,  $S = \{s_1, s_2, \dots, s_s\}$ , number of test genes,  $n \in \mathbb{Z}^+$ , a cut-off threshold,  $\tau \in \mathbb{Z}^+$ , and the number of allowed iterations,  $l \in \mathbb{Z}$

**Output:** Set of potential disease genes, denoted by  $A$

**Initialization:**

- Set  $i = 1$ ,  $A = \emptyset$ ,  $R = \{r_1, r_2, \dots, r_{n-s}\}$  – a set of randomly chosen genes, training set  $TR = O$ , test set  $TS = S \cup R$

**For**  $i \leq l$  **do**

1. Run a gene prioritization suite using the training set  $TR$ , test set  $TS$ , and  $m$  similarity criteria
2. Run  $k$  aggregation methods on the p-values produced in Step 1, and denote the resulting rankings by  $\sigma_1, \dots, \sigma_k$
3. Let  $B = \{\sigma_1(1), \dots, \sigma_1(\tau)\} \cup \dots \cup \{\sigma_k(1), \dots, \sigma_k(\tau)\}$
4.  $A \leftarrow A \cup B$ ;  $TR \leftarrow TR \cup B$ ;  $S \leftarrow S \setminus B$
5.  $TS \leftarrow S \cup R'$ ,  $R' =$  set of  $n - |S|$  randomly chosen genes
6.  $i \leftarrow i + 1$

**End**

**Return**  $A$

---

# Chapter 5

## RESULTS

We performed extensive cross-validation studies for eight diseases using both Endeavour- and ToppGene-generated p-values. Our results indicate that the similarity criteria that exhibit the strongest influence on the performance of the ToppGene and the Endeavour method are the PubMed and literature criteria, which award genes according to their citations in the disease related publications. In order to explore this issue further, we performed additional cross-validation studies for both ToppGene and Endeavour datasets to examine how exclusion of the literature criteria changes the performance of the two methods as well as our hybrid schemes. Our results reveal that HyDRA aggregation methods outperform Endeavour and ToppGene procedures for a majority of quality criteria, but they also highlight that each method offers unique advantages in prioritization for some specific diseases.

For gene discovery, we again used Endeavour and ToppGene p-values, and investigated three diseases – glioblastoma, meningioma and breast cancer – including all criteria available. We recommend as best practice a nested aggregation method, i.e., aggregating the aggregates of Endeavour, HyDRA, and ToppGene, coupled with iterative training set augmentation.

### 5.1 Cross-Validation

Cross-validation for HyDRA methods was performed on autism, breast cancer, colorectal cancer, endometriosis, ischaemic stroke, leukemia, lymphoma, and osteoarthritis. Tables 5.1 and 5.2 provide the summary of our results, pertaining to the average rank of one selected target gene. Table 5.1 illustrates that HyDRA methods offer optimal performance in 11 out of 16 tests when compared to ToppGene aggregates, and Table 5.2 illustrates that HyDRA outperforms Endeavour in 12 out of 16 cases. In the former case, the Weighted Hybrid Kendall method outperformed all other techniques.



Table 5.1: Cross-validation result of ToppGene and HyDRA methods for 8 diseases. Diseases without “\*” refer to results using all similarity categories in ToppGene. Diseases indexed by “\*” denote results which did not use the “Human Phenotype, Mouse Phenotype, Pubmed, Drug, Disease” similarity criteria in ToppGene. The scores describing the best average rank are bolded and shaded.

Disease	No. onset genes	ToppGene	Lovasz-Bregman	Hybrid Borda	Hybrid Kendall
Autism	40	7.275	11.2	9.75	<b>6.85</b>
Autism*	40	21.675	25.4	<b>19.775</b>	21.65
Breast Cancer	10	4.6	7.1	12	<b>2.5</b>
Breast Cancer*	10	<b>6.9</b>	17.8	8.1	7.1
Colorectal Cancer	20	7.3	<b>5.2</b>	7.85	8.7
Colorectal Cancer*	20	13.35	<b>9.5</b>	19.6	12.5
Endometriosis	43	<b>6.46</b>	8.63	10.63	7.74
Endometriosis*	43	<b>9.53</b>	9.76	15.84	9.7
Ischaemic Stroke	44	<b>5.61</b>	7.25	9.25	6.05
Ischaemic Stroke*	44	8.43	<b>7.5</b>	12.8	8.7
Leukemia	10	<b>5.5</b>	12	6.6	10.2
Leukemia*	10	20.8	22.8	24.3	<b>20.5</b>
Lymphoma	42	3.74	6.45	9.26	<b>2.93</b>
Lymphoma*	42	7.71	9.55	10.71	<b>6.76</b>
Osteoarthritis	41	6.44	6.51	13.54	<b>5.41</b>
Osteoarthritis*	41	8.73	8.32	14.1	<b>8.02</b>

Table 5.2: Cross-validation result of Endeavour and HyDRA methods for 8 diseases. Diseases without “\*” refer to results using all similarity categories in Endeavour. The indexing by “\*” corresponds to exclusion of similarity criteria “Precalculated-Ouzounis, Precalculated-Prospectr, Text” on Endeavour data. The scores describing the best average rank are bolded and shaded.

Disease	No. onset genes	Endeavour	Lovasz-Bregman	Hybrid Borda	Hybrid Kendall
Autism	40	17.96	19.3	17.78	<b>16.9</b>
Autism*	40	23.35	24.5	24.38	<b>21.78</b>
Breast Cancer	10	14.4	15	<b>12.5</b>	15.7
Breast Cancer*	10	16.6	<b>12.8</b>	15.5	17.8
Colorectal Cancer	20	8.55	8.65	<b>7.8</b>	8.1
Colorectal Cancer*	20	9.75	10.65	<b>9.55</b>	11.2
Endometriosis	43	5.3	6.37	<b>4.81</b>	5.65
Endometriosis*	43	<b>6.12</b>	7.63	6.86	6.6
Ischaemic Stroke	44	6.18	7.3	7.07	<b>6.09</b>
Ischaemic Stroke*	44	<b>7.95</b>	9.66	9.86	8.86
Leukemia	10	13.7	14.8	<b>7.1</b>	12.1
Leukemia*	10	19.5	19.9	<b>16.6</b>	21.3
Lymphoma	42	9.57	10.69	9	<b>8.81</b>
Lymphoma*	42	12.52	12.9	13.67	<b>11.67</b>
Osteoarthritis	41	<b>5.56</b>	6.32	7.46	6.29
Osteoarthritis*	41	<b>6.41</b>	7.41	6.51	7.22

Note that for all eight diseases, we performed two tests, in one of which we excluded those similarity criteria that contain strong prior information about disease genes, such as the “Disease” and “PubMed” category. Table 5.1 and Table 5.2 demonstrate the significant differences in average ranks of the target genes when literature information is excluded, suggesting that ToppGene and Endeavour both significantly benefit from this prior onset gene information when ranking the target genes. For true “discovery” methods one would usually not have such priors available, and results using these similarity criteria have to be treated with caution.

Another means for evaluating the performance of HyDRA algorithms compared to that of ToppGene and Endeavour is to examine the receiver operating characteristic (ROC) curves of the techniques. In this setting, we follow the same approach as used by both ToppGene and Endeavour. Sensitivity is defined as the frequency of tests in which prospect genes were ranked above a particular threshold position, and specificity as the percentage of prospect genes ranked below this threshold. As an example, a sensitivity/specificity pair of values 90/77 indicates that the presumably correct disease gene was ranked among the top-scoring  $100 - 77 = 23\%$  of the genes in 90% of the prioritization tests. The ROCs plot the dependence between sensitivity and the reflected specificity, and the area under the curve (AUC) represents another useful performance measure. The higher the AUC and specificity, the better the performance of the method. Figure 5.1 and Figure 5.2 are the ROC curves for cross validations comparing HyDRA with Endeavour and ToppGene, respectively. Endeavour reported 90/74 sensitivity/specificity values for their chosen set of test and training genes, as well as an AUC score of 0.866. Similarly, ToppGene reported 90/77 sensitivity/specificity values, and an AUC score of 0.916 for their tests of interest. Our specificity/sensitivity and AUC values are listed in Table 5.3 and Table 5.4, with best AUC and Sensitivity/Specificity values shaded in gray. Note that although the AUC values appear close in all cases, the HyDRA methods have very low overall computational complexity.

Table 5.3: AUC and Sensitivity/Specificity values for ToppGene and HyDRA rankings, pertaining to diseases listed in Table 5.1 using all criteria.

	ToppGene	Lovasz-Bregman	Hybrid Borda	Hybrid Kendall
AUC	0.951	0.93	0.911	0.947
Sensitivity/Specificity	90/84	90/75	90/75	90/84

Table 5.4: AUC and Sensitivity/Specificity values for Endeavour and HyDRA rankings, pertaining to diseases listed in Table 5.2 using all criteria.

	Endeavour	Lovasz-Bregman	Hybrid Borda	Hybrid Kendall
AUC	0.908	0.899	0.918	0.91
Sensitivity/Specificity	90/69	90/63	90/79	90/72

## 5.2 Gene Discovery

The genetic factors behind Glioblastoma, the most common and aggressive primary brain tumor, are still unknown. We study this disease, as well as meningioma and breast cancer, in the gene discovery phase. Our choice is governed by the fact that few publications are available pointing towards the causes of this form of brain cancer, and by the fact that it is widely believed that the genetic base of this disease is related to the genetic base of the Von Hippel-Lindau (VHL), Li-Fraumeni (LF), and Turcot syndromes (TS), neurofibromatosis (N), and tuberous sclerosis (TS) [32]. Furthermore, recent findings [33] indicate that brain cancers and breast cancers share a common line of mutations in the family of immunoglobulin GM genes, and that the human cytomegalovirus (HCMV) puts patients at risk of both brain and breast cancer.

Consequently, we used genes documented to be involved in glioblastoma as training genes for three discovery tests. In the first test, for the suspect genes we selected a subset of 15 genes known to be implicated in the VHL, LF, TS, N and TS syndromes. We subsequently ran Algorithm1 with  $l = 3$ ,  $s = 15$ ,  $n = 100$ ,  $\tau = 3$ . In the second test, we selected 18 genes known to be involved in breast cancer as suspect genes for glioblastoma, and ran Algorithm1 with  $l = 3$ ,  $s = 18$ ,  $n = 100$ ,  $\tau = 3$ . Finally, we performed the same analysis on suspect genes known to be involved in meningiomas, by

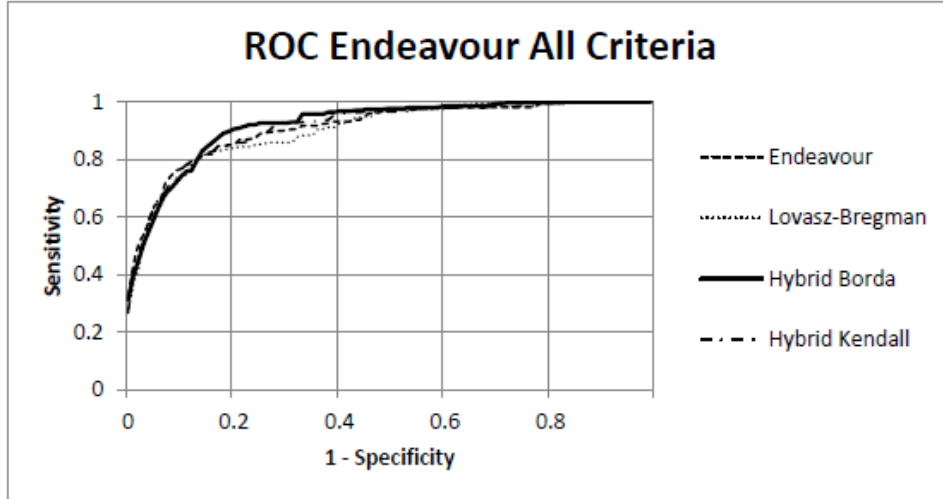


Figure 5.1: Cross-validation results: ROC curves for disease listed in Table 5.2 using all criteria and Endeavour data.

setting the parameters of iterative HyDRA gene discovery to  $l = 3$ ,  $s = 19$ ,  $n = 100$ ,  $\tau = 3$ . The results are shown in Table 5.5. Note that in our algorithmic investigation, we used  $l = 3$  (i.e., top-three) ranked genes, since this parameter choice offered a good trade-off between the size of the union of the top-ranked genes and the accuracy of the genes produced by the HyDRA discovery methods. The number of suspect genes was governed by the size of the available pool in OMIM/GAD and was targeted to be roughly 20% of the size of the test set. Such a percentage is deemed to be sufficiently high to allow for meaningful discovery, yet sufficiently low to prevent routine gene identification.

Table 5.5 reveals a number of results currently not known from the literature. The genes KRAS and CDH1, both implicated in breast cancer and meningioma, as well as CCND1 involved in meningioma<sup>1</sup> appear to be highly similar to genes implicated with glioblastoma. KRAS is a gene encoding for the K-Ras protein that is involved in regulating cell division, and hence an obvious candidate for being implicated in cancer. On the other hand, CDH1 is responsible for the production of the E-cadherin protein, whose function is to aid in cell adhesion and to regulate transmission of chemical signals within cells, and control cell maturation. E-cadherin also often acts as a tumor suppressor protein. GeneCards reveals that the CCND1 gene is implicated in

<sup>1</sup>As well as in colorectal cancer.

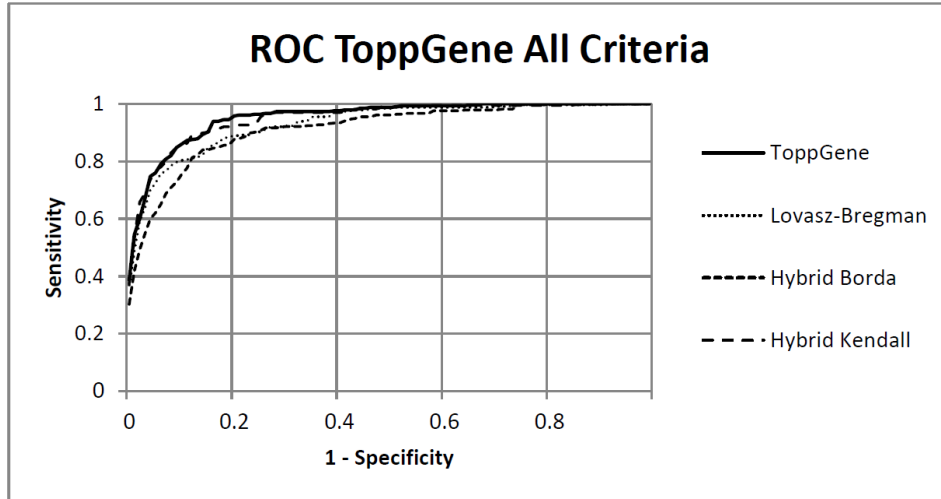


Figure 5.2: Cross-validation results: ROC curves for disease listed in Table 5.1 using all criteria and ToppGene data.

altering cell cycle progression, and is mutated in a variety of tumors. Its role in glioma tumorigenesis appears to be well documented [34], but surprisingly, neither KRAS nor CDH1 nor CCND1 are listed in the OMIM/GAD database as potential glioblastoma genes.

Another interesting finding involves genes ranked among the top three candidates, but not identified as “suspect” genes. For instance, according to GeneBank, GSTM2 regulates an individual’s susceptibility to carcinogens and toxins and may suggest glioblastoma being in part caused by toxic and other environmental conditions; KAAG1 appears to be implicated with kidney tumors, while TP73 belongs to the p53 family of transcription factors and is known to be involved in neuroblastoma.

Table 5.5: The union of top 3 ranked genes from ToppGene, Endeavour, and HyDRA methods for the three suspect gene discovery sets, with the “suspect” genes in bold. In all cases, the training genes are genes implicated in Glioblastoma. The “Disease” category indicates from which family of diseases the test genes were drawn. VHL = Von Hippel-Lindau, LF = Li-Fraumeni, TS = Turcot Syndromes, N = Neurofibromatosis, and TS = Tuberos Sclerosis.

Test Disease	Iteration 1	Iteration 2
Breast Cancer	<b>AKT1</b> , ATM, <b>BRIP1</b> , <b>CDH1</b> , <b>CHEK2</b> , GSTM2, KAAG1, <b>RAD51</b> , TP73	<b>BARD1</b> , CASP7, ITGA4, <b>KRAS</b> , <b>PALB2</b> , <b>PHB</b> , SMAD7, UMOD
VHL, LF, TS, N, TS	<b>CCND1</b> , CD28, CD74, <b>CDK4</b> , <b>CHEK2</b> , <b>MLH1</b> , <b>MSH2</b> , <b>MSH6</b> , NBPF4, <b>PMS2</b> , PRNT, <b>TSC2</b>	ALCAM, <b>APC</b> , <b>MRC1</b> , NCL, <b>NF1</b> , <b>NF2</b> , SNCA, TAF7, TOPBP1, <b>TSC1</b> , <b>VHL</b>
Meningioma	<b>CCND1</b> , <b>HLA-DQB1</b> , KLF6, <b>KRAS</b> , <b>TGFB1</b> , <b>TGFBR2</b> , <b>XRCC5</b>	BAGE, <b>BAP1</b> , CAV1, CD4, <b>CDH1</b> , <b>NF2</b> , <b>PDGFB</b> , PSMC2, <b>RFC1</b> , SAMD9L, SERPING1, <b>SMARCB1</b>

# Chapter 6

## CONCLUSION

We start by discussing the results in Table 5.1 and Table 5.2. The first observation is that the Lovász-Bregman method performs worse than any other aggregation method. This finding may be attributed to the fact that the p-values have a large span, and small values may be “masked” by larger ones. Scaling all p-values may be a means to improve the performance of this technique, but how exactly to accomplish this task remains a question.

In almost all cases, except for leukemia and lymphoma, the average rankings produced by ToppGene and the weighted Kendall distance appear to be almost identical. But *average* values may be misleading, as individual rankings of genes may vary substantially between the methods. It is for this reason that we recommend merging lists generated by different methods as best aggregation practice. Another important observation is that HyDRA methods have significantly lower computational complexity than ToppGene and, especially, Endeavour, and hence scale well for large datasets.

Another finding is the fact that the good performance of ToppGene and all other methods largely depends on including prior literature on the genes into the aggregation process. We observed situations where the rank of an element dropped by roughly 90 positions when this prior was not available. This implies that for gene discovery, it is risky to rely on any single method, and it is again good practice to merge top-ranked entries generated by different methods. Finally, it is not clear how to optimally choose the number of training genes for a given set of test genes, or vice versa. Choosing more training genes may appear to be beneficial at first glance, but it creates a more diverse pool of candidates for which some similarity criteria will inevitably fail to identify the right genes. In this case, we recommend using the weighted Kendall to eliminate outliers, and in addition, we recommend the use of a fairly large TvB scaling parameter.



## REFERENCES

- [1] G. M. Cooper and R. E. Hausman, *The Cell*. Sinauer Associates Sunderland, 2000.
- [2] M. Kim, F. Farnoud, and O. Milenkovic, “Hydra: gene prioritization via hybrid distance-score rank aggregation,” *Bioinformatics*, vol. 31, no. 7, pp. 1034–1043, 2015.
- [3] E. A. Adie, R. R. Adams, K. L. Evans, D. J. Porteous, and B. S. Pickard, “Suspects: enabling fast and effective prioritization of positional candidates,” *Bioinformatics*, vol. 22, no. 6, pp. 773–774, 2006.
- [4] L. R. Cardon and J. I. Bell, “Association study designs for complex diseases,” *Nature Reviews Genetics*, vol. 2, no. 2, pp. 91–99, 2001.
- [5] N. Risch, K. Merikangas et al., “The future of genetic studies of complex human diseases,” *Science-AAAS-Weekly Paper Edition*, vol. 273, no. 5281, pp. 1516–1517, 1996.
- [6] N. Tiffin, E. Adie, F. Turner, H. G. Brunner, M. A. van Driel, M. Oti, N. Lopez-Bigas, C. Ouzounis, C. Perez-Iratxeta, M. A. Andrade-Navarro et al., “Computational disease gene identification: a concert of methods prioritizes type 2 diabetes and obesity candidate genes,” *Nucleic Acids Research*, vol. 34, no. 10, pp. 3067–3081, 2006.
- [7] E. Guney, J. Garcia-Garcia, and B. Oliva, “Guildify: a web server for phenotypic characterization of genes through biological data integration and network-based prioritization algorithms,” *Bioinformatics*, p. btu092, 2014.
- [8] D. Warde-Farley, S. L. Donaldson, O. Comes, K. Zuberi, R. Badrawi, P. Chao, M. Franz, C. Grouios, F. Kazi, C. T. Lopes et al., “The gene-mania prediction server: biological network integration for gene prioritization and predicting gene function,” *Nucleic Acids Research*, vol. 38, no. 2, pp. W214–W220, 2010.
- [9] C. Perez-Iratxeta, P. Bork, and M. A. Andrade, “Association of genes to genetically inherited diseases using data mining,” *Nature Genetics*, vol. 31, no. 3, pp. 316–319, 2002.

- [10] F. S. Turner, D. R. Clutterbuck, C. A. Semple et al., “Pocus: mining genomic sequence annotation to predict disease genes,” *Genome Biology*, vol. 4, no. 11, pp. R75–R75, 2003.
- [11] S. Aerts, D. Lambrechts, S. Maity, P. Van Loo, B. Coessens, F. De Smet, L.-C. Tranchevent, B. De Moor, P. Marynen, B. Hassan et al., “Gene prioritization through genomic data fusion,” *Nature Biotechnology*, vol. 24, no. 5, pp. 537–544, 2006.
- [12] T. De Bie, L.-C. Tranchevent, L. M. Van Oeffelen, and Y. Moreau, “Kernel-based data fusion for gene prioritization,” *Bioinformatics*, vol. 23, no. 13, pp. i125–i132, 2007.
- [13] J. Chen, E. E. Bardes, B. J. Aronow, and A. G. Jegga, “Toppgene suite for gene list enrichment analysis and candidate gene prioritization,” *Nucleic Acids Research*, vol. 37, no. 2, pp. W305–W311, 2009.
- [14] T. Kacprowski, N. T. Doncheva, and M. Albrecht, “Networkprioritizer: a versatile tool for network-based prioritization of candidate disease genes or other molecules,” *Bioinformatics*, p. bt164, 2013.
- [15] P. C. Fishburn, “Arrow’s impossibility theorem: Concise proof and infinite voters,” *Journal of Economic Theory*, vol. 2, no. 1, pp. 103–106, 1970.
- [16] J. G. Kemeny, “Mathematics without numbers,” *Daedalus*, vol. 88, no. 4, pp. 577–591, 1959.
- [17] M. Popescu, J. M. Keller, and J. A. Mitchell, “Fuzzy measures on the gene ontology for gene product similarity,” *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, vol. 3, no. 3, pp. 263–274, 2006.
- [18] M. G. Kendall, “A new measure of rank correlation,” *Biometrika*, pp. 81–93, 1938.
- [19] M. G. Kendall, *Rank correlation methods*. Oxford: Griffin, 1948.
- [20] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar, “Rank aggregation methods for the web,” in *Proceedings of the 10th International Conference on World Wide Web*. ACM, 2001, pp. 613–622.
- [21] E. Thanassoulis, *Introduction to the Theory and Application of Data Envelopment Analysis*. Springer, 2001.
- [22] F. Farnoud, B. Touri, and O. Milenkovic, “Nonuniform vote aggregation algorithms,” in *Signal Processing and Communications (SPCOM), 2012 International Conference on*. IEEE, 2012, pp. 1–5.

- [23] J. J. Bartholdi III, C. A. Tovey, and M. A. Trick, “The computational difficulty of manipulating an election,” *Social Choice and Welfare*, vol. 6, no. 3, pp. 227–241, 1989.
- [24] H. W. Kuhn, “The Hungarian method for the assignment problem,” *Naval Research Logistics Quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [25] F. Farnoud and O. Milenkovic, “An axiomatic approach to constructing distances for rank comparison and aggregation,” *IEEE Transactions on Information Theory*, vol. 60, no. 10, pp. 6417–6439, 2014.
- [26] F. Raisali, F. F. Hassanzadeh, and O. Milenkovic, “Weighted rank aggregation via relaxed integer programming,” in *Information Theory Proceedings (ISIT), 2013 IEEE International Symposium on*. IEEE, 2013, pp. 2765–2769.
- [27] “The Hungarian algorithm.” 2006. [Online]. Available: <http://www.mathworks.com/matlabcentral/fileexchange/11609-hungarian-algorithm>
- [28] R. Iyer and J. Bilmes, “The Lovász-Bregman divergence and connections to rank aggregation, clustering, and web ranking: Extended version,” *Uncertainty in Artificial Intelligence*, 2013.
- [29] A. Hamosh, A. F. Scott, J. S. Amberger, C. A. Bocchini, and V. A. McKusick, “Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders,” *Nucleic Acids Research*, vol. 33, no. 1, pp. D514–D517, 2005.
- [30] K. G. Becker, K. C. Barnes, T. J. Bright, and S. A. Wang, “The genetic association database,” *Nature Genetics*, vol. 36, no. 5, pp. 431–432, 2004.
- [31] M. Safran, I. Solomon, O. Shmueli, M. Lapidot, S. Shen-Orr, A. Adato, U. Ben-Dor, N. Esterman, N. Rosen, I. Peter et al., “Genecards 2002: towards a complete, object-oriented, human gene compendium,” *Bioinformatics*, vol. 18, no. 11, pp. 1542–1543, 2002.
- [32] A. P. Kyritsis, M. L. Bondy, J. S. Rao, and C. Sioka, “Inherited predisposition to glioma,” *Neuro-oncology*, vol. 12, no. 1, pp. 104–113, 2009.
- [33] J. P. Pandey, “Immunoglobulin gm genes, cytomegalovirus immunoevasion, and the risk of glioma, neuroblastoma, and breast cancer,” *Frontiers in Oncology*, vol. 4, p. 236, 2014.
- [34] R. Büschges, R. G. Weber, B. Actor, P. Lichter, V. P. Collins, and G. Reifenberger, “Amplification and expression of cyclin d genes (ccnd1 ccnd2 and ccnd3) in human malignant gliomas,” *Brain Pathology*, vol. 9, no. 3, pp. 435–442, 1999.