

A Decision-theory Approach to Interpretable Set Analysis for High-dimensional Data

Simina M. Boca, Hector Corrada Bravo, Jeffrey T. Leek
Johns Hopkins University
Bloomberg School of Public Health

Giovanni Parmigiani
Dana Farber Cancer Institute
Harvard School of Public Health

Abstract

A ubiquitous and challenging problem in the analysis of high-dimensional data is to identify pre-defined sets of features that are enriched for features that show an association of interest. Here we propose a new decision theory approach to gene set analysis that addresses the most common problems in the set-level analysis of high-dimensional data including: (1) overlapping set annotations, (2) non-comparable p-values across sets of varying sizes, and (3) competitive tests. Instead of calculating an enrichment p-value for each set, we search for estimators, consisting of unions of disjoint sets, that minimize a weighted average of the number of false discoveries and missed discoveries. We show that this minimization problem can be written in terms of a sum of the posterior probabilities that the data for each feature comes from the alternative distribution. We also define a new false discovery rate called the “atomic false discovery rate” for non-overlapping sets and show that the loss function which weights the number of false discoveries and the number of missed discoveries has a simple analytic solution based on thresholding the atomic false discovery rate at a fixed level. We apply our approach to a microarray study of breast cancer and show that our decision theory approach improves the analysis of intersecting gene sets.

1 Introduction

Modern scientific studies often measure a large number of variables on each sample. These variables are usually measurements of certain physical properties, or “features.” In genomics, the features may be genes and the actual variables gene expression measurements; tens or hundreds of thousands of gene expression measurements are taken on only a few subjects. Similarly, in brain imaging studies, tens of thousands of voxel intensities are measured on a small number of study participants. In spatial epidemiology, estimates of disease prevalence are sometimes calculated for a large number of locations, with a small number of individuals in each location.

One common approach to high-dimensional data analysis is to perform inference on a set of features indexed by $\mathcal{M} = \{1, 2, \dots, M\}$ marginally, one at a time. The data for the m th feature is an $N \times 1$ matrix, which is compared to a single $N \times 1$ outcome vector Y . Feature level statistics are functions of the data X_m and the outcome Y . The null hypothesis for a given feature is that it is not strongly associated with the outcome of interest. Thus, each feature can be seen as coming either from the null or an alternative distribution. Shrinkage methods (Newton, Kendziorski, Richmond, Blattner, and Tsui (2001), Cui, Hwang, Qiu, Blades,

and Churchill (2005)) and empirical Bayes (Efron and Tibshirani (2002), Newton and Kendziorski (2003), Gottardo, Pannucci, Kuske, and Brettin (2003), Newton, Noueir, Sarkar, and Ahlquist (2004)) use the full set of feature data $X = \{X_m\}_{m=1}^M$, however the focus is still on marginal analysis of individual features.

In biology, there is also a keen interest in set-level inference (Tavazoie, Hughes, Campbell, Cho, and Church (1999), Mirnics, Middleton, Marquez, Lewis, and Levitt (2000), Bouton and Pevsner (2002)). Set analyses of high-dimensional data have been developed to combine information across features, to infer associations between sets and outcomes or phenotypes. Sets are defined by previous experimentally verified or postulated relationships between features. Usually, a fixed number of sets K is defined in advance. Each of the sets $\mathcal{S} = \{S_1, S_2, \dots, S_K\}$ is a subset of \mathcal{M} . Some sets are based on the spatial or geographic arrangement of features and are naturally disjoint: For instance, a voxel may get mapped to exactly one brain region, and a geographical location will get mapped to exactly one county, but in general features may belong to more than one set. For example, one gene may participate in more than one metabolic pathway.

Set-level inference can be useful in a number of ways. It can increase statistical power, by helping detect correlated changes in features within the same set, which are too subtle to be detected in a marginal analysis (Mootha, Lindgren, Eriksson, Subramanian, Sihag, Lehar, Puigserver, Carlsson, Ridderstrale, Laurila et al. (2003)). This can be due to coregulation of gene expression within pathways, voxel intensities within brain regions, or disease prevalences in locations within broader geographic regions, such as counties. Another possibility is that the signal is at the level of the sets rather than the features. For example, the same pathway may be altered in different tumors, but the specific genes which are altered may differ between tumors (Parsons, Jones, Zhang, Lin, Leary, Angenendt, Mankoo, Carter, Siu et al. (2008)). Inferences at the level of sets are also sought to increase interpretability: For instance, hundreds of genes may be up-regulated in a specific phase of cell division in bacteria, though most of those genes may belong to a small number of known pathways. Knowing that only a few pathways are involved is much more informative than knowing that hundreds of genes show varying levels of expression.

Here we propose a general decision-theoretic approach for set-level analysis. We introduce atoms, which are non-overlapping sets derived from the original set annotation. We propose to deal with only these atoms, to avoid the interpretive and statistical difficulties associated with standard gene set methods discussed later. For statistical inference we focus on the expected number of false and missed discoveries, rather than using p-values. These can be interpreted in relation to the estimated prevalence of features from the alternative within a given atom. Our approach combines these two key ideas into a decision theory framework where we

derive estimators consisting of unions of atoms that minimize a weighted combination of the EFD and EMD. These estimators are functions of the joint posterior distribution of feature-specific parameters and do not require a separate statistical treatment of feature-level and set-level analysis.

The results we present here can be summarized as follows.

1. Section 3 introduces atoms.
2. Section 4 outlines the decision-theory framework we are using. In particular, in Theorem 3 we show that the two components of the main loss function we can consider, which represent the posterior expected number of false discoveries and the posterior expected number of missed discoveries, may be written in terms of only the marginal feature-level posterior probabilities.
3. Section 5 provides alternative loss functions which may be used, including loss functions with regularization penalties and a loss function which uses the ratio of false discoveries and missed discoveries.
4. Section 6 introduces the concept of atomic FDR and provides algorithms for finding the Bayes estimators for the loss functions described above. For most of the loss functions considered, an analytic closed form is provided.
5. Section 7 outlines an empirical Bayes approach for estimating the feature-level posterior probabilities and applies this approach to simulation and real gene expression analyses.

2 Relationship to previous approaches

A standard approach for feature-level analysis is to perform a statistical test for each feature to assess the level of association with the outcome of interest. For example, in genomics a test may be used to decide whether a gene is differentially expressed between two conditions. This can be a t-test or a F-test, or some variation on them (Storey and Tibshirani (2003a), Baldi and Long (2001), Cui et al. (2005)). Due to the large number of tests performed, a multiple comparison adjustment is needed. The goal is often to control the false discovery rate (FDR) (as in Benjamini and Hochberg (1995), Efron and Tibshirani (2002), Storey (2002), Storey and Tibshirani (2003b), Storey (2003)). Alternatively, the posterior probabilities that features are from the null or alternative distributions, or from mixture components of interest, can be estimated using Bayes (Parmigiani, Garrett, Anbazhagan, and Gabrielson (2002), Do, M'uller, and Tang (2005)) or Empirical Bayes methods (L'onnstedt and Speed (2002), Efron and Tibshirani (2002), Newton and Kendziorski (2003), Gottardo et al. (2003), Newton et al. (2004)). These methods generally borrow

information across genes, but the inference is still performed marginally, one feature at a time, and no information from outside the experiment is considered.

Most set-level inference methods combine the feature-level statistics into set-level statistics, then perform some hypothesis test for each set. In genomics, one approach is to individually score genes, then separate them into categories based on these scores; the simplest scenario is to choose a single cutoff and declare the genes above the cutoff “differentially expressed” and the genes below the cutoff “non-differentially expressed.” For each set a contingency table is then constructed, which cross-classifies the genes by whether or not they are in the set and by their category and a classical statistical test is performed, such as the hypergeometric test (Tavazoie et al. (1999)). Alternatively, a test can be performed to compare the distributions of statistics of the genes within a set to those of all the genes (Mirnics et al. (2000)). A method with a somewhat different flavor, *gene-set enrichment analysis* (GSEA), is described in Subramanian, Tamayo, Mootha, Mukherjee, Ebert, Gillette, Paulovich, Pomeroy, Golub, Lander et al. (2005) (a preliminary version is employed in Mootha et al. (2003)). It starts with a list of ordered gene scores. It then calculates an enrichment score for each gene-set by using a signed variation of the Kolmogorov-Smirnov statistic. A p-value is calculated for the gene-set by permuting the phenotype labels. A normalized enrichment score is then computed, taking into account the gene-set size, and a FDR-control method is employed. A similar method is implemented in the *limma* package in *R*. (Smyth (2004)), which uses the average of the scores of genes in the gene-set instead of the signed Kolmogorov-Smirnov statistic and obtains p-values by permuting the genes instead of the phenotype labels. A speedier implementation is available which only uses the ranks of the genes, as opposed to the scores. In this case, the Wilcoxon test is used to obtain a p-value, and permutations are not necessary.

Despite the appeal of set level analyses, there remain four significant difficulties that must be overcome when analyzing these data. (1) One problem is that the set annotations often overlap, which can cause issues with the interpretability of the results. Furthermore, p-value methods also require a multiple testing adjustment, which can be greatly complicated by the correlation that results from set overlaps. (2) Another problem results from the confusing interpretation p-values often have, as they represent the probability of getting a result as or more extreme than the one observed, and thus do not directly convey information about the fraction of features in a set which are from the alternative distribution, which we call the fraction of alternatives. In particular, p-values have different interpretations for different sample sizes, and generally set annotations cannot be expected to produce sets which all consists of the same number of features. This can lead to sets with very different fractions of alternatives having similar p-values. (3) A third problem results from the common use of “competitive tests” (Goeman and Buhlmann

(2007)) in many set-level inference procedures, which pit a set against its complement. This results in the “zero sum problem,” where more significant features in one set often lead to higher p-values in other sets. (4) Lastly, most approaches perform separate analyses at the feature and set level, so uncertainty remaining from the feature-level analysis is ignored at the set level.

To expand on the second point, a large number, if not the majority, of methods for set-level inference rely on calculating a p-value for each set. However, p-values do not actually give a direct indication of how much signal is present in a given set. Several common methods for set-level inference in genomics which rely on calculating p-values are reviewed and critiqued in the influential paper of Geman and Buhlmann (2007). We revisit some of their points here, namely the issues of gene sampling and competitive tests, and show how our approach bypasses these problems. Their proposed method involves considering instead a null hypothesis that no genes in a gene-set represent true signal. However, this does not appear to be the null hypothesis that most scientists would be interested in. We consider that it is more informative, in the case of set-level inference, to move away from a hypothesis framework to an estimation framework.

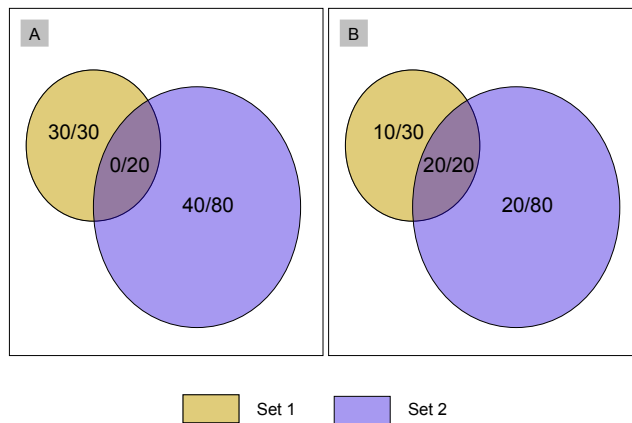
Gene sampling methods use permutations of the gene, as opposed to the sample, labels to obtain a p-value. This means that a theoretical replication of the experiment would involve considering a different set of genes which are subjected to the same measurements on the same subjects. In particular, this means that the sample size is given by the number of genes, as opposed to the number of subjects. One of the practical problems which results from this approach is that p-values have different interpretations for different sample sizes, with larger sample sizes often leading to smaller p-values, even if there is less signal present. Thus, two sets of different sizes with different fractions of alternatives may end up with very similar p-values.

With regard to the third point, competitive tests consider the null hypothesis that features in a given set are at most as often differentially expressed as features in the complement of the set. This results in a “zero sum problem,” where a larger fraction of alternatives in one set may lead to higher p-values in other sets. Methods based on the Fisher exact test, the chi-squared test, the Wilcoxon rank test, and the z-test fall into this class. Consider, for example, the case where the signal in a dataset is represented by 10% of the features. Then take sets consisting of 10, 50, and 100 features, none of which have any signal. The larger sets will then have larger p-values, because as the size of the set increases, the size of its complement will decrease, and therefore the proportion of signal in the complement will increase.

3 Atoms

Overlapping sets of features may result in erroneous statistical inferences by inducing correlations between sets. For example, in the case of methods which are based on calculating a p-value for each set, multiple testing adjustments can be complicated by correlations between the hypotheses. Overlap can also cause issues with interpretability, e.g. if a large set in \mathcal{S} gets a poor score, but a smaller set which shares a large fraction of its genes with the large gene-set gets a good score, it is unclear where this difference comes from. See Figure 1 for an example in terms of gene-sets and their fraction of differentially expressed genes: In both (A) and (B), the smaller set (set 1) has 50 genes, 30 of which are differentially expressed, the larger set (set 2) has 100 genes, 40 of which are differentially expressed, and their intersection has 20 genes. However, in panel (A), none of the genes in the intersection is differentially expressed, whereas in (B), they are all differentially expressed. The percent of differentially expressed genes in the difference between sets 2 and 1 varied greatly between the two cases (50% and 25%, respectively).

Figure 1: Set 1 has 60% of genes differentially expressed (30/50), while set 2 has 40% of genes differentially expressed (40/100). However, in (A) there are no differentially expressed genes common to sets 1 and 2, and although set 2 has a lower percentage of differentially expressed genes compared to set 1, the percentage of differentially expressed genes which are in set 2 but not in set 1 is higher than the percentage of differentially expressed genes which are in set 2 (50% compared to 40%). In (B), the percentage of differentially expressed genes which are in set 2 but not in set 1 (25%) is lower than the percentage of differentially expressed genes in set 2 (40%).



Here we consider only non-overlapping sets, which we call *atoms*. By considering only non-overlapping sets we avoid the difficulties in analysis and interpretation from overlap. We let \mathcal{A} be the set of atoms which we obtain from the sets, \mathcal{S} where $\mathcal{A} = \{A_1, \dots, A_L\}$ has L different non-overlapping sets. Each element of \mathcal{A} is a subset of an element in \mathcal{S} and the intersection of any two elements in \mathcal{A} is empty. One way to obtain atoms is to divide the sets into their smallest non-overlapping subunits, so that the collection of atoms \mathcal{A} is defined as the set of minimal cardinality which has the following properties:

1. Given any set $S \in \mathcal{S}$, there exists $\mathcal{J} \in \{1, \dots, L\}$ such that $S = \cup_{i \in \mathcal{J}} A_i$.
2. Given any atoms A_i and A_j in \mathcal{A} with $i \neq j$, $A_i \cap A_j = \emptyset$.

In the examples shown in Figure 1 (A) and (B), this is a very natural step to take, as we can simply consider three atoms, consisting of the intersection between sets 1 and 2, the difference between set 1 and set 2, and the difference between set 2 and set 1. In particular, in Figure 1 (A), there appear to be two disjoint sets with large fractions of alternatives, whereas in Figure 1 (B), there seems to be a common signal between sets 1 and 2. Hence, it is more informative to split the overlapping sets. An alternative way of obtaining atoms is by clustering the features using a dissimilarity measure which depends on the sets in \mathcal{S} which are shared between features (Boca et al., in preparation).

4 Decision theory framework

We consider a decision theory framework to integrate the feature-level and set-level inference. While the concept is general, we detail it here in terms of a loss function that is linear in two components, one relating to false discoveries and the other to missed discoveries. Within this context we study conditions under which the posterior expected loss can be written in terms of feature-level posterior probabilities.

The set of features from the alternative distribution is denoted by τ , which is a subset of \mathcal{M} . For example, in the case where the features are gene expression measurements, τ represents the genes which are differentially expressed between the two conditions. If we wanted to estimate τ based on only the feature-level information, we would simply take the estimate to be the set of all features whose statistics in X are within certain limits. For instance, if we performed a two-sample t-test for a differential gene expression problem, we could estimate τ by taking all the genes with p-values below some threshold, or with p-values below some threshold and effect size (in this case, fold change) above another threshold.

Estimating τ based on just the feature-level data would ignore the scientific background which is represented by the atom-level information. Thus, we seek to

estimate τ using the unions of atoms in \mathcal{A} . We denote by \mathcal{U} the set of all possible unions of atoms in \mathcal{A} . We look for the set in \mathcal{U} which maximizes the overlap with τ by using a relevant loss function. Thus, we are in the situation of providing a scientifically meaningful estimate of τ by solving a constrained estimation problem where we only have elements in \mathcal{U} as possible estimators.

The discrepancy between τ and a candidate estimator U in \mathcal{U} can be represented by a loss function. We consider the following general class of loss functions, which depend on a discrepancy function d and a fixed constant $w \in [0, 1]$:

$$L(\tau, U) = (1 - w) \sum_{m \in U \setminus \tau} d(m, \tau) + w \sum_{m \in \tau \setminus U} d(m, U) \quad (1)$$

for all $U \in \mathcal{U}$. We note that $U \setminus \tau$ represents the set of *false discoveries*, while $\tau \setminus U$ represents the set of *missed discoveries* if we were to estimate τ by U . Thus, the loss function is linear in two components, the first one measuring how close features which are false discoveries are to the set of features from the alternative distribution (τ), the second one measuring how close features which are missed discoveries are to the candidate estimator (U).

We consider general discrepancy measures d between features in \mathcal{M} which satisfy:

$$\begin{aligned} d : \mathcal{M} \times \mathcal{M} &\rightarrow [0, \infty) \\ d(m, m) &= 0 \text{ for any } m \in \mathcal{M} \\ d(m_1, m_2) &= d(m_2, m_1) \text{ for any } m_1, m_2 \in \mathcal{M} \end{aligned}$$

We could use for example the single-linkage (nearest neighbor) function to define the discrepancy between a feature and a set $d(m, A) = \min\{d(m, m_0) : m_0 \in A\}$.

Based on the loss function L , we get the following posterior expected loss:

$$\mathcal{L}(U) = \sum_{\tau \in 2^{\mathcal{M}}} L(\tau, U) * P(\tau|X, Y)$$

where $2^{\mathcal{M}}$ is the power set of \mathcal{M} and $P(\tau|X, Y)$ is the posterior probability of set τ exactly representing the set of features which are from the alternative distribution. In particular, $P(\tau|X, Y)$ can be seen as a discrete probability distribution with support $2^{\mathcal{M}}$. It fully reflects uncertainty from the feature-level modeling and dependencies between features.

We use the following notation for the posterior probabilities, for simplicity:

$$p_\tau = P(\tau|X, Y)$$

Thus, the posterior expected loss may be written as:

$$\begin{aligned}\mathcal{L}(U) &= (1-w) \sum_{\tau \in 2^{\mathcal{M}}} \sum_{n \in U \setminus \tau} d(m, \tau) p_{\tau} + w \sum_{\tau \in 2^{\mathcal{M}}} \sum_{n \in \tau \setminus U} d(m, U) p_{\tau} \\ &= (1-w) E_{\tau|X,Y} \left\{ \sum_{m \in U \setminus \tau} d(m, \tau) \right\} + w E_{\tau|X,Y} \left\{ \sum_{m \in \tau \setminus U} d(m, U) \right\}\end{aligned}$$

The two components may be interpreted as the posterior expected value of the sum of distances from each false discovery ($m \in U \setminus \tau$) to the set of features from the alternative distribution (τ) and the posterior expected value of the sum of distances from each missed discovery ($m \in \tau \setminus U$) to the set of all discoveries (U).

The posterior expected loss is written in terms of posterior set-level probabilities. In order to obtain a Bayes estimator, we would need to minimize it, which could be extremely complicated from a modeling point of view, because we would need to model the joint distribution of all features. It would also be extremely computationally intensive, as the number of unions of atoms is 2^L , where we recall that L is the total number of atoms.

It would be much easier to estimate the posterior expected loss for a particular set and to find the union of atoms which minimizes it if we could write it in terms of the posterior probabilities of individual features being from the alternative distribution. We call these probabilities *marginal feature-level posterior probabilities*. They can be estimated in a relatively straight-forward manner, by building a probability model and using a fully Bayesian framework or an Empirical Bayes (EB) framework.

We introduce the following notation for the *marginal posterior probability for a set U* , which is the sum of all posterior probabilities of sets which include U :

$$p_U^* = \sum_{\tau \in 2^{\mathcal{M}}, U \subset \tau} p_{\tau}$$

It represents the posterior probability that the set U is included in the set of all features from the alternative distribution. For specific cases we can simply write out the features in the set, e.g. $p_{12} = p_{21} = p_{\{1,2\}}$. The marginal feature-level posterior probabilities are a specific case of this, where the set represents a single feature, i.e.:

$$p_m^* = \sum_{m \in \tau} p_{\tau}$$

We prove a lemma which simplifies the form of the two components of the posterior expected loss function. Note, in particular, that $E_{\tau|X,Y} \left\{ \sum_{m \in \tau \setminus U} d(m, U) \right\}$ can be written as a linear function of marginal feature-level posterior probabilities.

Lemma 1. *Under the loss function described in equation (1), in the case of a general discrepancy measure d and single linkage, the following simplified forms of $E_{\tau|X,Y}\{\sum_{m \in U \setminus \tau} d(m, \tau)\}$ and $E_{\tau|X,Y}\{\sum_{m \in \tau \setminus U} d(m, U)\}$ are obtained:*

$$\begin{aligned} E_{\tau|X,Y}\left\{\sum_{m \in U \setminus \tau} d(m, \tau)\right\} &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in U \setminus \tau} d(m, \tau) p_{\tau} = \sum_{m \in U} \sum_{\tau \in 2^{\mathcal{M}}, m \notin \tau} d(m, \tau) p_{\tau} \\ E_{\tau|X,Y}\left\{\sum_{m \in \tau \setminus U} d(m, U)\right\} &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in \tau \setminus U} d(m, U) p_{\tau} = \sum_{m \notin U} d(m, U) p_m^* \end{aligned}$$

Proof.

$$\begin{aligned} \sum_{m \in U} \sum_{\tau \in 2^{\mathcal{M}}, m \notin \tau} d(m, \tau) p_{\tau} &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in U \setminus \tau} d(m, \tau) p_{\tau} = \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in U \setminus \tau} d(m, \tau) p_{\tau} \\ \sum_{m \notin U} d(m, U) p_m^* &= \sum_{m \notin U} d(m, U) \sum_{\tau \in 2^{\mathcal{M}}, m \in \tau} p_{\tau} = \sum_{m \in \tau \setminus U, \tau \in 2^{\mathcal{M}}} d(m, U) p_{\tau} \\ &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in \tau \setminus U} d(m, U) p_{\tau} \end{aligned}$$

□

We now show that $E_{\tau|X,Y}\{\sum_{m \in U \setminus \tau} d(m, \tau)\}$ cannot be written as an affine function of marginal feature-level posterior probabilities for a general discrepancy measure d which takes into account how far or close features are to each other, in the case where the single linkage property holds. Thus, in general, to calculate the posterior expected loss, we need to model the joint distribution of all the features. This also leads to much more complex computations.

Lemma 2. *Under the loss function described in equation (1), in the case of a general discrepancy measure d and single linkage, $E_{\tau|X,Y}\{\sum_{m \in U \setminus \tau} d(m, \tau)\}$ cannot be written as an affine function of marginal feature-level posterior probabilities. Therefore $\mathcal{L}(U)$ also cannot be written as an affine function of marginal feature-level posterior probabilities.*

Proof. We show that we cannot write $E_{\tau|X,Y}\{\sum_{m \in U \setminus \tau} d(m, \tau)\}$ as an affine function of the marginal feature-level posterior probabilities.

Step 1. We first show that, for any proper subset $v \subsetneq 2^{\mathcal{M}}$, setting any affine of the posterior probabilities of the sets in v equal to 0 forces all the coefficients to be 0, i.e.:

Denote by elements in ν by $\tau_{l_1}, \dots, \tau_{l_{|\nu|}}$. We will show that setting any affine function of these elements to 0 implies that all the coefficients are 0. Thus, we have:

$$\sum_{\tau \in \nu} a_{\tau} p_{\tau} + b = 0 \quad (2)$$

We note that if $\nu = \{\tau_{l_1}, \dots, \tau_{l_{|\nu|}}\} \subsetneq 2^{\mathcal{M}}$, then $p_{\tau_{l_1}} + \dots + p_{\tau_{l_{|\nu|}}} \leq 1$ and $p_{\tau_{l_1}} \geq 0, \dots, p_{\tau_{l_{|\nu|}}} \geq 0$. Plugging in $p_{\tau_{l_1}} = 1, p_{\tau_{l_2}} = \dots = p_{\tau_{l_{|\nu|}}} = 0$, followed by $p_{\tau_{l_1}} = \frac{1}{2}, p_{\tau_{l_2}} = \dots = p_{\tau_{l_{|\nu|}}} = 0$, and solving the resulting system of equations in $a_{\tau_{l_1}}$ and b results in $a_{\tau_{l_1}} = b = 0$. From here on, plugging in only one non-zero probability for each $\tau \in \nu$ in turn will result in $a_{\tau_{l_1}} = \dots = a_{\tau_{l_{|\nu|}}} = 0$.

Step 2. We now apply the result in *Step 1* to show that $E_{\tau|X,Y}\{\sum_{m \in U \setminus \tau} d(m, \tau)\}$ can in general not be written as an affine function of the marginal feature-level posterior probabilities. We note that we have:

$$E_{\tau|X,Y}\left\{\sum_{m \in U \setminus \tau} d(m, \tau)\right\}(U) = \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in U \setminus \tau} d(m, \tau) p_{\tau} = \sum_{\tau \in 2^{\mathcal{M}}, \tau \neq U} \left\{\sum_{m \in U \setminus \tau} d(m, \tau)\right\} p_{\tau} \quad (3)$$

since $d(m, U) = 0$ if $m \in \tau$. Using a simple transformation, we note that showing that $E_{\tau|X,Y}\{\sum_{m \in U \setminus \tau} d(m, \tau)\}$ we need to show that we can find a_m and b such that:

$$\begin{aligned} E_{\tau|X,Y}\left\{\sum_{m \in U \setminus \tau} d(m, \tau)\right\}(U) &= \sum_{m \in \mathcal{M}} a_m (1 - p_m^*) + b \\ &= \sum_{m \in \mathcal{M}} a_m \sum_{\tau \in 2^{\mathcal{M}}, m \notin \tau} p_{\tau} + b \\ &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \notin \mathcal{M} \setminus \tau} a_m p_{\tau} + b \\ &= \sum_{\tau \in 2^{\mathcal{M}}, \tau \neq \mathcal{M}} \left\{\sum_{m \in \mathcal{M} \setminus \tau} a_m\right\} p_{\tau} + b \end{aligned} \quad (4)$$

The coefficients a_m and b are more accurately written as $a_m(U)$ and $b(U)$, but we use the simpler notation here. Setting the expressions in 3 and 4 equal to each other, we get:

$$\sum_{\tau \in 2^{\mathcal{M}}, \tau \neq U} \left\{\sum_{m \in U \setminus \tau} d(m, \tau)\right\} p_{\tau} = \sum_{\tau \in 2^{\mathcal{M}}, \tau \neq \mathcal{M}} \left\{\sum_{m \in \mathcal{M} \setminus \tau} a_m\right\} p_{\tau} + b$$

We may now take $\nu = 2^{\mathcal{M}} \setminus \mathcal{M}$, which is a proper subset of $2^{\mathcal{M}}$. Using the result in *Step 1*, we get:

$$\sum_{m \in U \setminus \tau} a_m = \sum_{m \in U \setminus \tau} d(m, \tau),$$

all the other coefficients being 0. Now consider cycling through all the sets τ such that $U \setminus \tau$ consists of a single element. We thus obtain:

$$d(m, \tau) = a_m \text{ for all } m \in U \setminus \tau$$

regardless of how many elements there are in $\tau \setminus U$ and how far away they are from the elements in $U \setminus \tau$.

To illustrate this last portion of the proof, consider $\mathcal{M} = \{1, 2, 3\}$ and $U = \{1, 2\}$. Then:

$$\begin{aligned} \tau = \{2\} & \Rightarrow d(1, \{2\}) = a_1 \\ \tau = \{2, 3\} & \Rightarrow d(1, \{2, 3\}) = a_1 \end{aligned}$$

Given our use of the single-linkage property, $d(1, \{2, 3\}) = \min\{d(1, \{2\}), d(1, \{3\})\}$. So if $d(1, \{3\}) \geq d(1, \{2\})$, then $d(1, \{2\}) = d(1, \{3\})$, which means that the discrepancy measure d does not take into account how far or close features are to each other. \square

To make the analysis tractable, we consider the scenario where d does not take into account how far or close features are to each other. In this case, the discrepancy is equivalent to the following 0 – 1 function which takes as inputs two features m_1 and m_2 :

$$d(m_1, m_2) = 1(m_1 \neq m_2). \quad (5)$$

This is equivalent to saying that the discrepancy between two features is 0 if and only if the features are one and the same, otherwise it is 1. We note that in this case the loss function is reduced to:

$$\begin{aligned} L(\tau, U) &= (1 - w) \sum_{m \in U \setminus \tau} d(m, \tau) + w \sum_{m \in \tau \setminus U} d(m, U) \\ &= (1 - w) * |U \setminus \tau| + w * |\tau \setminus U| \\ &= (1 - w) * \text{Number of false discoveries} + w * \text{Number of missed discoveries} \end{aligned}$$

When considering the posterior expected loss function, the two components become the expected number of false discoveries and the expected number of missed discoveries, which we will denote by $EFD(U)$ and $EMD(U)$:

$$\mathcal{L}(U) = (1 - w)EFD(U) + wEMD(U)$$

Similar loss functions have been used by Storey (2003) and Müller, Parmigiani, Robert, and Rousseau (2004), though not in set-level inference.

In Theorem 3, we simplify the form of both $EFD(U)$ and $EMD(U)$. In particular, both can be written as affine functions of the marginal feature-level posterior probabilities. This yields major benefits in terms of both modelling and computation.

Theorem 3. *Under the loss function described in equation (1) and the single linkage property, $EFD(U)$ and $EMD(U)$ can both be written as affine functions of the marginal feature-level posterior probabilities. $EFD(U)$ can be written as the sum of posterior probabilities that the features in U are from the null distribution and $EMD(U)$ can be written as the sum of the posterior probabilities that the features which are not in U are from the alternative distribution:*

$$\begin{aligned} EFD(U) &= \sum_{m \in U} (1 - p_m^*) = |U| - \sum_{m \in U} p_m^* \\ EMD(U) &= \sum_{m \notin U} p_m^* \end{aligned}$$

Proof. Using the lemma above:

$$\begin{aligned} EFD(U) &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in U \setminus \tau} p_\tau = \sum_{m \in U} \sum_{\tau \in 2^{\mathcal{M}}, m \notin \tau} p_\tau = \sum_{m \in U} (1 - p_m^*) = |U| - \sum_{m \in U} p_m^* \\ EMD(U) &= \sum_{\tau \in 2^{\mathcal{M}}} \sum_{m \in \tau \setminus U} p_\tau = \sum_{m \notin U} p_m^* \end{aligned}$$

□

From this section onward, we only consider this 0 – 1 discrepancy measure defined in (5) and the corresponding loss function, due to the many advantages it has over other discrepancy measures. We note that, although it has a relatively simple form, it still results in a flexible modelling framework.

5 Alternative loss functions

We consider some variations on the loss function corresponding to the 0 – 1 discrepancy measure which was introduced earlier, by looking at what happens when linear constraints are introduced. In order to penalize for a large number of features in the Bayes estimator, we consider the following constrained optimization problem:

Minimize $L(\tau, U)$ with the constraint that $|U| < \rho$, for some $\rho > 0$. This is equivalent to minimizing the loss function:

$$L_f^\lambda(\tau, U) = (1 - w) * |U \setminus \tau| + w * |\tau \setminus U| + \lambda |U| \text{ for some } \lambda > 0$$

(see for instance Gill, Murray, and Wright (1981)).

Similarly, one may be interested in penalizing for a large number of atoms in the Bayes estimator. For example, in a spatial epidemiology example, one may desire to target an intervention to a fixed number of locations due to issues of cost or work force. In such a case, we would consider the following constrained optimization problem:

Minimize $L(\tau, U)$ with the constraint that $|J| < \eta$, for some $\eta > 0$, where J is the number of atoms in U . This is equivalent to minimizing the loss function:

$$L_a^\xi(\tau, U) = (1 - w) * |U \setminus \tau| + w * |\tau \setminus U| + \xi |J| \text{ for some } \xi > 0.$$

We remark that we could also penalize unions of atoms with a small number of genes, respectively atoms, by simply changing the signs of λ and ξ . We could also have a loss function which incorporates penalties on both the number of features and the number of atoms.

Alternatively, loss functions which employ the ratio of missed discoveries and false discoveries, instead of the number of missed discoveries and false discoveries, may also be used. Thus, we can consider:

$$\begin{aligned} L_r(\tau, U) &= (1 - w) * \frac{|U \setminus \tau|}{|U|} + w * \frac{|\tau \setminus U|}{M - |U|} \\ &= (1 - w) * \text{Ratio of false discoveries} + w * \text{Ratio of missed discoveries} \end{aligned}$$

6 Atomic FDR

In general, finding the Bayes estimator for the loss functions based on the 0 – 1 discrepancy measure can be very computationally challenging, even after the simplification in Theorem 3, since the number of sets in \mathcal{U} which the posterior expected loss needs to be minimized over is 2^L , where L is the total number of atoms. We find an analytic solution for obtaining the Bayes estimator for the loss function which weights the number of false discoveries and missed discoveries, as well as the loss functions with regularization penalties, described in Section 5. This result, established in Theorem 5, shows that the Bayes estimator for the loss L is found by choosing those atoms $A_l \in \mathcal{A}$ with EFD less than or equal to w . This EFD can be thought of as a “realized” (Müller, Parmigiani, and Rice (2007)) or Bayesian (?)

false discovery rate for atom A_l , which we denote by FDR_l , and call the realized *atom false discovery rate* (atomic FDR). Thus, the algorithm for finding the Bayes estimator corresponds to thresholding the atomic FDR at a fixed level determined by the parameter w . For large values of w the procedure allows more false positives, since the EFD is down-weighted in the loss function, while small values of w more strongly weight the EFD and restrict the resulting atomic false discovery rate. In the case where we penalize atoms with many genes, considering L_f^λ , the Bayes estimator is equivalent to thresholding the realized atomic FDR adjusted for a “background rate” of false discoveries λ .

We first prove a result which gives a convenient parametrization of the posterior expected loss using the original loss function:

Lemma 4. *The posterior expected loss $\mathcal{L}(U)$ which results from the 0 – 1 dissimilarity measure may be rewritten as:*

$$t' \{(1-w)n - q\} + w1'q$$

where t, q, n , and 1 are vectors of length $L = |\mathcal{A}|$. They may be written out as:

$$\begin{aligned} t &= (t_1, t_2, \dots, t_L) \\ q &= (q_1, q_2, \dots, q_L) \\ n &= (n_1, n_2, \dots, n_L) \\ 1 &= (1, 1, \dots, 1) \end{aligned}$$

where t_l is the indicator of whether atom A_l is part of U , n_l is the number of features in atom A_l ($n_l = |A_l|$), and q_l is the sum of the marginal posterior probabilities of the genes in atom A_l , i.e. $q_l = \sum_{m \in A_l} p_m^*$.

Proof.

$$\begin{aligned} \mathcal{L}(U) &= (1-w) \sum_{m \in U} (1-p_m^*) + w \sum_{m \notin U} p_m^* \\ &= (1-w) \sum_{A_l \in U} (n_l - q_l) + w \sum_{n_l \notin U} q_l \\ &= (1-w)t'(n - q) + w(1-t)'q \\ &= (1-w)t'n - (1-w)t'q + w1'q - wt'q \\ &= t' \{(1-w)n - q\} + w1'q \end{aligned}$$

□

The posterior expected loss functions \mathcal{L}_f^λ and \mathcal{L}_a^ξ which correspond to the loss functions $L_f^\lambda(\tau, U)$ and $L_a^\xi(\tau, U)$ (from Section 5) have similar parametrizations, which are also linear in t :

$$\mathcal{L}_f^\lambda(U) = t' \{(1-w)n - q\} + w1'q + \lambda t'n \quad (6)$$

$$\mathcal{L}_a^\xi(U) = t' \{(1-w)n - q\} + w1'q + \xi 1'n \quad (7)$$

We provide a straightforward algorithm for obtaining the Bayes estimator for any loss function for which the posterior expected loss can be parametrized as a affine function of t , including L , L_f^λ , and L_a^ξ , for a fixed w between 0 and 1. This is because any affine function of t , $h(t) = t'\mathfrak{J} + \mathfrak{T}$, $t \in \{0, 1\}^J$, $\mathfrak{J} \in \mathbb{R}^J$, $\mathfrak{T} \in \mathbb{R}$ is minimized when $t_j = 1\{\mathfrak{J}_j \leq 0\}$, since $h(t)$ is equivalent to minimizing $t'\mathfrak{J}$. This is a linear function in each component t_j of t , and if we minimize it in each component we also minimize it overall. As a result, it is minimized by choosing to sum only over those components of \mathfrak{J} which are negative or zero.

Theorem 5. *For a fixed value of $w \in [0, 1]$, analytic solutions are obtained for the Bayes estimators for the losses L , L_f^λ , and L_a^ξ :*

The indicator t_l of whether atom A_l is in the Bayes estimator for the loss L is:

$$t_l = 1\left\{1 - \frac{q_l}{n_l} \leq w\right\}$$

The indicator t_l of whether atom A_l is in the Bayes estimator for the loss L_f^λ is:

$$t_l = 1\left\{1 - \left(\frac{q_l}{n_l} - \lambda\right) \leq w\right\}$$

The indicator t_l of whether atom A_l is in the Bayes estimator for the loss L_a^ξ is:

$$t_l = 1\left\{1 - \frac{q_l - \xi}{n_l} \leq w\right\}$$

For the loss function L_r from Section 5, which considers the ratio of false discoveries and missed discoveries, an analytic solution is not available. However, an approximate algorithmic solution is presented below:

We first consider the solution to the problem where we constrain the size of the Bayes estimator $|U| = n't$ to some size ρ . In this case, we get the following constrained linear binary problem:

$$\begin{aligned} \min_t \quad & t' \left\{ \frac{(1-w)}{\rho}(n-q) - \frac{w}{(M-\rho)}q \right\} \\ \text{s.t.} \quad & n't = \rho \end{aligned}$$

This is an instance of the well-known 0-1 knapsack problem (Garey and Johnson, 1979), which can be solved approximately by Dantzig's greedy algorithm. This uses a sorting strategy where atoms are sorted increasingly by the quantity

$$\frac{(1-w)}{\rho} \left(1 - \frac{q_i}{n_i} \right) - \frac{w}{(M-\rho)} \frac{q_i}{n_i}$$

and t_i is set to 1, in order, until $n't = \rho$. Note that when $\rho = M/2$, atoms are sorted according to $(1 - q_i/n_i)$ as in Theorem 5.

In principle, to solve the fractional problem, one can solve the 0-1 knapsack problem for each possible value of $\rho = |U|$ and select the best solution. Since ρ can range over a large number of possible values, we use a strategy based on the projected gradient of the fractional function at a given point to find a small number of estimator sizes $|U|$ to test.

In summary, the algorithm is as follows: (1) initialize ρ ; (2) find solution t_ρ by solving the 0-1 knapsack problem; (3) find point s as the minimizer of $L_r(t)$ along the linear-piecewise path $t_\rho - \alpha(\nabla_t L_r(t))_+$ where $(\nabla_t L_r(t))_+$ is the projected gradient at t_ρ (see Wright and Nocedal, 2006, for a nice illustration); (4) stop if $s'n = \rho$, otherwise set $\rho = s'n$ and repeat step (2). Step (3) is a univariate optimization problem which can be solved using any univariate numerical minimization technique (golden-section method, for instance).

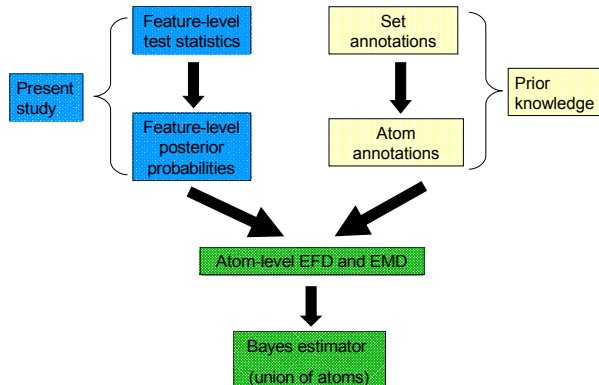
7 Applications

7.1 Estimation of posterior probabilities

We initially perform feature-level statistical inference, then translate it to the set level. The feature level statistics represent information from the present study or studies. We combine this information with prior knowledge, represented by annotations placing the features into sets. We only consider non-overlapping sets for

set-level inference. If the sets we start out with are overlapping, we use the annotations to obtain non-overlapping sets, which we call atoms (see section 3). We then calculate the expected missed discoveries and false discoveries for every atom, and use this to obtain a Bayes estimator corresponding to a particular loss function. This workflow is summarized in Figure 2. In our illustrations, we estimated the feature-

Figure 2: Workflow diagram. Note that we are combining two sources of information: The feature-level information from the present study, and the annotation information, which represents a distillation of prior scientific knowledge. They are combined to obtain the atom-level EFD (expected false discoveries) and EMD (expected missed discoveries), and finally, the Bayes estimator.



level posterior probabilities using an Empirical Bayes approach, following the lead of Efron and Tibshirani (2002) and Newton and Kendziorski (2003). In particular, we use the nonparametric approach detailed in Storey, Akey, and Kruglyak (2005):

1. Multiple sets of null statistics are obtained by using distributional assumptions. The observed statistics are denoted by F_m and the null statistics by F_{m0}^b , where m indexes the features, $1 \leq m \leq M$ and b the simulations under the null, $1 \leq b \leq B$.
2. The probability π_0 of a randomly selected feature being from the null distribution is estimated for a series of thresholds c by $\hat{\pi}_0(c) = \frac{\#\{F_m \leq c\}}{\#\{F_m^{0b} \leq c\}/B}$. The estimate $\hat{\pi}_0$ is then chosen by smoothing over $\hat{\pi}_0(c)$, employing the approach detailed in Storey and Tibshirani (2003b). In the cases where the number of features is less than 500, π_0 is conservatively assumed to be 1.

3. The density of the prior null distribution is denoted by g_0 and that of the alternative distribution by g_1 . Thus, the density of the F_m statistics is $g = \pi_0 g_0 + (1 - \pi_0) g_1$.
4. The null and observed statistics are concatenated, with the observed statistics (F_m) considered “successes” and the null statistics (F_{m0}^b) considered “failures.”
5. The ratio $g_0(F)/g(F)$ is estimated by logistic regression using the successes and failures defined in step 4, employing the approach in Anderson and Blair (1982), with a natural cubic spline using a fixed number of equally spaced knots, as in Green and Silverman (1994).
6. The posterior probability of a specific feature m being from the alternative distribution can be written as $1 - \pi_0 \frac{g_0(F_m)}{g(F_m)}$ and is thus easily estimated from steps 1 and 5.

We consider 20 simulations under the null to be sufficient and we take the number of equally spaced knots to be equal to the total number of features.

7.2 Simulations

We carried out simulations which compared our method to a method which relies on the Wilcoxon rank test and is available in the *limma* package in *R*. We calculated p-values for the Wilcoxon test, as well as q-values, which are adjusted p-values to control the FDR for independent hypothesis tests (Benjamini and Hochberg (1995)). In each simulation, the features which are from the alternative distribution are draws from a normal distribution with mean 1 and variance 1/15, while remaining features are draws from a normal distribution with mean 0 and variance 1/15.

We first revisited the scenario described in Figure 1. We performed 100 simulations with the 4 sets described in parts (A) and (B) of Figure 1. We give the results from the Wilcoxon rank test in Table 1. We note that, since Set 1A and Set 2A and Set 1B and Set 2B are overlapping, the results from this test are difficult to interpret, since it is unclear where the difference in p-values comes from. In fact, Set 1A and Set 2A have no features in common which are from the alternative distribution, while Set 1B and Set 2B have 20 such features in common. In Table 2 we present the results which use our method and are applied on the atoms we obtained from simply taking the intersections and differences of the original sets. The p-values for the individual atoms presented in Table 1 are much more interpretable, and the $1 - EFDR$ values even more so.

We performed another set of 100 simulations, with 2500 features, 5% of which were from the alternative distribution. These 125 features were distributed

Table 1: Summary from 100 simulations using the sets described in Figure 1, (A) and (B), and performing a Wilcoxon rank test. The mean and standard deviation of the p-values and q-values are calculated over the 100 runs. Set 1A and Set 2A have a fraction of alternatives of 0.6, while Set 1B and Set 2B have a fraction of alternatives of 0.4. Since the sets are overlapping, the test gives no information about where the large difference in p-values come from.

	Set	Fract. of alt.	mean p-value	sd p-value	mean q-value	sd q-value
1	1A	0.6	0.045	0.057	0.109	0.115
2	2A	0.4	0.887	0.099	0.96	0.049
3	1B	0.6	0.05	0.054	0.119	0.115
4	2B	0.4	0.883	0.089	0.968	0.04

Table 2: Summary from 100 simulations using the sets described in Figure 1, (A) and (B), and employing our method with the atoms obtained from the intersections and differences of the original sets. The mean and standard deviation of the estimated 1-EFDR and the mean and standard deviation of the p-values and q-values are calculated over the 100 runs. We note that the results are much more interpretable than those in Table 1.

	Atom	Fract. of alt.	mean 1 - EFDR	sd 1 - EFDR	mean p-value	sd p-value	mean q-value	sd q-value
1	$1A \setminus 2A$	1	0.916	0.038	<0.001	<0.001	<0.001	<0.001
2	$1A \cap 2A$	0	0.024	0.03	0.999	0.001	1	<0.001
3	$2A \setminus 1A$	0.5	0.467	0.019	0.294	0.166	0.572	0.301
4	$1B \setminus 2B$	0.33	0.318	0.028	0.864	0.112	0.993	0.033
5	$1B \cap 2B$	1	0.914	0.044	<0.001	<0.001	<0.001	<0.001
6	$2B \setminus 1B$	0.25	0.245	0.018	1	<0.001	1	<0.001

among 5 atoms of size 50, which had different fractions of alternatives, of 0.9, 0.7, 0.5, 0.3, and 0.1. We also considered an atom of size 50 with no features from the alternative distribution, for comparison. The remaining features were not placed in any atoms. The mean posterior probability over the 100 simulation runs

was estimated to be 0.068, with a standard deviation of 0.021, which is quite close to the true value of 0.05. In Table 3 we present the results which compare our method to the Wilcoxon rank test. We note that our method provides much more interpretable results. The Wilcoxon test gives an average q-value of 0.011 for a fraction of alternatives of 0.3, with a standard deviation of 0.03, which means that for many of the simulation runs this set would not be considered significant. This is despite the fact that a fraction of alternatives 0.3 is much higher than the background fraction of 0.05. This effect is even more pronounced for a fraction of alternatives of 0.1, where the mean q-value is 0.407.

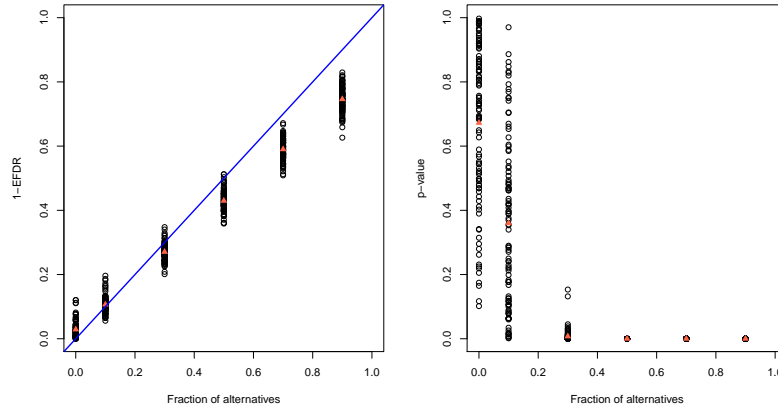
Table 3: Summary from 100 simulations with 6 atoms of size 50. There are 2500 features total, the 2200 which are not distributed among the atoms all being from the null distribution. Thus, the overall percentage of features which are from the alternative distribution is 5%. The mean and standard deviation of the estimated 1-EFDR, as well as the mean and standard deviation of the p-values and q-values from the Wilcoxon rank test, are calculated over the 100 runs.

	Fract. of alt.	mean 1-EFDR	sd 1-EFDR	mean p-value	sd p-value	mean q-value	sd q-value
1	0.9	0.746	0.039	<0.001	<0.001	<0.001	<0.001
2	0.7	0.59	0.034	<0.001	<0.001	<0.001	<0.001
3	0.5	0.43	0.032	<0.001	<0.001	<0.001	<0.001
4	0.3	0.271	0.026	0.007	0.021	0.011	0.03
5	0.1	0.107	0.027	0.36	0.263	0.407	0.285
6	0	0.029	0.027	0.672	0.247	0.686	0.238

Plots of $1 - EFDR$ and the q-values versus the true fractions of alternatives for the simulations presented in Table 3 are shown in Figure 3. In the ideal scenario, the mean value of $1 - EFDR$ would be close to the true fraction of alternatives. Given that the posterior probabilities are estimated to be between 0 and 1, we expect a slight anti-conservative bias for the sets with low fractions of alternatives and a conservative bias for the sets with high fractions of alternatives, which is what we see in the plot in the left panel. The plot in the right panel shows that the p-values have a very wide spread for the low fractions of alternatives, but there is nearly no spread for the higher fractions of alternatives, highlighting the difference in interpretation between estimated fraction of alternatives and significance tests.

We also explored the effect of the atom size on both our method and the Wilcoxon rank test method. We considered 2500 features again, 5% of which are

Figure 3: True fractions of alternatives for the simulated datasets described in Table 3. The red triangles represent the mean values over the 100 simulation runs for $1 - EFDR$ (left panel) and for the p-values (right panel). The blue line in the left panel represents the ideal scenario, where $1 - EFDR$ perfectly estimates the fraction of alternatives.



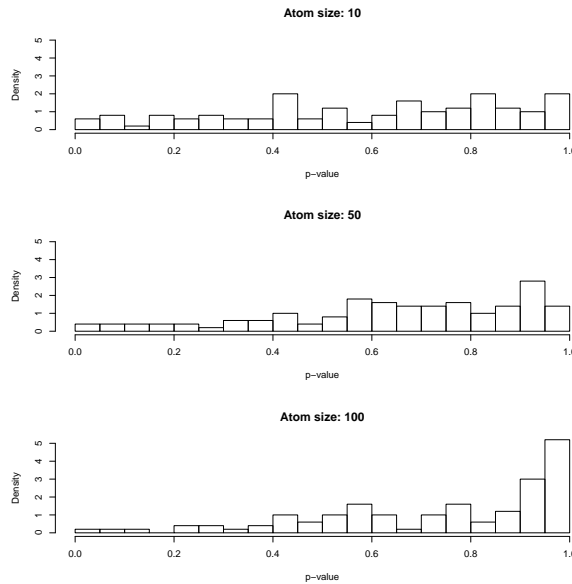
from the alternative distribution. We took 3 atoms, each having a fraction of alternatives of 0.5, but different sizes: 10, 50, and 100 features. The results are displayed in Table 4. For our method, the mean estimate over the 100 runs is similar across the different set sizes. As expected, for the Wilcoxon rank test method, the p-values and q-values for the atom of size 10 is much higher than for the other two atoms. As noted in Section 2 and confirmed in this simulation, the p-values in set inference methods which employ feature sampling have different interpretations for different set sizes. Thus, the Wilcoxon rank test is not particularly interpretable in this case. Our decision theoretic, estimation-oriented framework provides a much clearer interpretation.

Table 4: Summary from 100 simulations with 3 atoms having fractions of alternatives of 0.5, but different set sizes. The mean and standard deviation of the estimated $1 - EFDR$, as well as the mean and standard deviation of the p-values and q-values from the Wilcoxon rank test, are calculated over the 100 runs.

	Size	mean 1-EFDR	sd 1-EFDR	mean p-value	sd p-value	mean q-value	sd q-value
1	10	0.44	0.059	0.019	0.031	0.019	0.031
2	50	0.436	0.034	<0.001	<0.001	<0.001	<0.001
3	100	0.438	0.026	<0.001	<0.001	<0.001	<0.001

We also considered the impact of the atom size for atoms which have only null features. Once again, we used 2500 features total, 5% of which were from the alternative distribution. As in the previous example, we considered atoms of sizes 10, 50, and 100. Whereas our method gives an estimated mean EFDR between 0.961 and 0.985 for each of the three atoms, with standard errors smaller than 0.04, the Wilcoxon rank test results in p-values which are increasingly skewed towards 1 as the atom size increases. This is due to the feature-sampling and competitive nature of the Wilcoxon test, as noted in Section 2: for atoms with no features from the alternative, as the atom size increases, the fraction of features from the alternative distribution in the complement of the atom increases. In Figure 4, we show histograms exhibiting this behavior for the different sets, over 100 runs.

Figure 4: Histograms of the p-values obtained from the Wilcoxon rank test for atoms of sizes 10, 50, and 100 which have no features from the alternative distribution. Note as the set size increases, the histograms are increasingly skewed towards 1. The mean p-value for the atom of size 10 is 0.591, while for the atom of size 50 it is 0.652, and for the atom of size 100 it is 0.737.



We also compared the Bayes estimators resulting from the loss function L which weights the number of false discoveries and missed discoveries to those resulting from the other loss functions we introduced in Section 5, namely L_f^λ , L_a^ξ , and L_r . We compared the results on 100 simulations with 2250 features, 10% of

which were from the alternative distribution. We considered 8 atoms, 4 of size 50 and 4 of size 100. For each of the two sizes considered, atoms had fractions of alternatives of 0, 0.1, 0.5, or 0.9. The results are presented in Table 5. We used $\lambda = 0.20$ and $\xi = 5$. The most frequently selected Bayes estimators from the 100 runs were compared to the ideal scenario, where the features from the alternative and null distributions are given posterior probabilities of 1, respectively 0. We note that the interpretation of w for the loss function L_r is different from that for the other three loss functions. In general, the estimators which were commonly chosen in the simulation runs were subsets of the one in the ideal scenario.

Table 5: Bayes estimators obtained from different loss functions over 100 runs. Atom names use the fraction of alternatives and the set size; for example atom 0.5 - 50 has a fraction of alternatives 0.5 and 50 features. The ideal scenario (in bold), gives posterior probabilities of 1 and 0 to the features from the alternative, respectively from the null distribution. It is compared to the simulation results. In parentheses are listed the number of simulations in which a particular union of atoms is the Bayes estimator (only if it appeared in at least 20 out of 100 runs.)

w	L	L_f^λ	L_a^ξ	L_r
0.25	0.9 - 50, 0.9 - 100	Empty set	0.9 - 50, 0.9 - 100	0.9 - 50, 0.9 - 100
	0.9 - 50, 0.9 - 100 (79)	Empty set (100)	Empty set (68)	0.9 - 50, 0.9 - 100 (61)
			0.9 - 100 (27)	0.9 - 50 (37)
0.5	0.9 - 50, 0.5 - 50, 0.9 - 100, 0.5 - 100	0.9 - 50, 0.9 - 100	0.9 - 50, 0.9 - 100	0.9 - 50, 0.9 - 100
	0.9 - 50, 0.9 - 100 (87)	0.9 - 50, 0.9 - 100 (97)	0.9 - 50, 0.9 - 100 (100)	0.9 - 50, 0.9 - 100 (84)
0.67	0.9 - 50, 0.5 - 50, 0.9 - 100, 0.5 - 100	0.9 - 50, 0.9 - 100	0.9 - 50, 0.5 - 50, 0.9 - 100, 0.5 - 100	0.9 - 50, 0.5 - 50, 0.1 - 50, 0.9 - 100, 0.5 - 100, 0.1 - 100
	0.9 - 50, 0.5 - 50, 0.9 - 100, 0.5 - 100 (100)	0.9 - 50, 0.9 - 100 (99)	0.9 - 50, 0.5 - 50, 0.9 - 100, 0.5 - 100 (79)	0.9 - 50, 0.5 - 50, 0.1 - 50, 0.9 - 100, 0.5 - 100, 0.1 - 100 (52)
				0.9 - 50, 0.9 - 100 (36)

7.3 Data analysis

We present an analysis of a dataset from Subramanian et al. (2005), which compares mRNA expression profiles from lymphoblastoid cell lines of 15 males and 17 females. This data was originally analyzed via the GSEA method, and was later analyzed with a different method, which used t-tests, in Irizarry, Wang, Zhou, and Speed (2009). The gene-sets used represented chromosomal regions. For illustration, we excluded 40 of the original 212 sets, in order to obtain nonoverlapping atoms. We compared the methods in Subramanian et al. (2005) and Irizarry et al. (2009), as well as the Wilcoxon rank test, to our method. The results from the top ten sets using our method are presented in Table 6.

Table 6: Comparison of our method to the GSEA, t-test, and Wilcoxon rank test methods on a dataset from Subramanian et al. (2005). We present the results for the top 10 sets using our method.

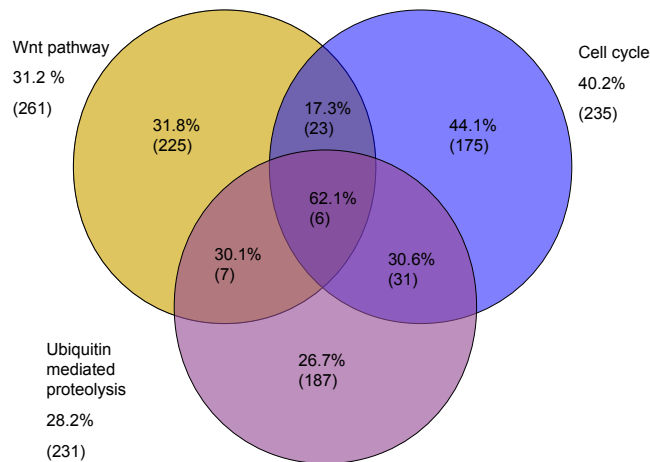
	Set	EFDR	rank test p-values	rank test q-value	GSEA p-value	GSEA q-value	t-test p-value	t-test q-value
1	chrYq11	0.621	<0.001	0.014	<0.001	<0.001	<0.001	<0.001
2	chr4q22	0.956	0.154	0.995	0.195	0.996	0.176	0.92
3	chr2q14	0.958	0.867	0.995	0.411	0.996	0.302	0.92
4	chrXp11	0.959	0.204	0.995	0.149	0.996	0.058	0.641
5	chrXq26	0.963	0.999	0.999	0.677	0.996	0.79	0.984
6	chr14q21	0.965	0.939	0.995	0.902	0.996	0.742	0.984
7	chr10q21	0.966	0.939	0.995	0.54	0.996	0.358	0.92
8	chr7p22	0.966	0.92	0.995	0.891	0.996	0.986	0.992
9	chr5q33	0.966	0.872	0.995	0.263	0.996	0.328	0.92
10	chr12p12	0.967	0.103	0.995	0.636	0.996	0.893	0.984

The overall estimated expected fraction of false discoveries is 0.976. The set which has the lowest EFDR with our method is the only set which had a chromosomal region on the Y chromosome, and it also ranked first in terms of p-values and q-values with the other two methods. Given that the primary difference between the two groups is gender, it can serve as a proof of principle. We note that our method is much more interpretable than methods which rely on p-values or q-values: While the estimate of the fraction of alternatives in the set chrYq11 is substantially higher than both the overall estimate and the estimate for the set ranked second, it might

not be considered extremely high in absolute terms. Therefore, providing the actual estimate and allowing direct comparisons appears to be much more useful than trying to understand the difference between a q-value of nearly 0 and a q-value of over 0.99.

We also analyzed a dataset from Sotiriou, Wirapati, Loi, Harris, Fox, Smeds, Nordgren, Farmer, Praz, Haibe-Kains et al. (2006), consisting of expression microarrays from breast tumors. We looked at a subset of untreated tumors, considering the differential expression between ER-positive and ER-negative samples. 10 of the samples were ER-negative, 53 were ER-positive. Using the KEGG annotations for the Wnt pathway, cell cycle, and ubiquitin mediated proteolysis, we saw that the strongest signal is found in the atom represented by their intersection, as opposed to one of the 3 sets. Thus, as seen in Figure 5, using atoms provides us with much more informative results.

Figure 5: The estimated percentage of genes from the alternative $((1 - EFDR) * 100)$ for 3 sets and the atoms resulting from them, using data from Sotiriou et al. (2006). The numbers in parentheses represent the number of genes in the sets or atoms. Note that the atom created from the intersection of the 3 sets has an estimated percentage of genes from the alternative of 62.1%, while the percentages for the original sets were all lower than 41%. Thus, the important information in this case seem to lay in the intersection of the three sets, which can be interpreted an interaction.



8 Discussion

We introduced a general approach for set-level inference for high-dimensional data, which casts the problem in a decision-theoretic framework and focuses on estimation rather than testing. Set-level inference is an area of increasing interest in many areas of science, because of the necessity of combining quantitative feature-level data with annotations resulting from alternative sources of information. Our method introduces the concept that set-level inference is best performed for disjoint sets (atoms), in order to obtain increased scientific clarity and interpretability. We discuss in detail an implementation that focuses on quantifying the differences between sets based on the expected number of false discoveries (EFD) and the expected number of missed discoveries (EMD). These have a clear interpretation and provide information about the question of greatest interest, which relates to quantifying the fraction of alternatives in each set.

Our approach introduces a new paradigm in set-level inference. Most present methods are based on performing a hypothesis test for each set. The p-values thus obtained are fed into a set-level analysis which in turn requires a multiple testing adjustment. The statistical properties of this overall strategy are difficult to interpret. We provide a rigorous unified framework for feature and set-level analysis. Our estimates have clearly defined optimality properties and are scientifically interpretable.

We show that the loss function defined as the weighted sum of false discoveries and missed discoveries for any union of atoms, can be reduced to a form which depends only on the marginal feature-level posterior probabilities. These probabilities can easily be estimated using existing Empirical Bayes methods. This simplification enables us to obtain an easy algorithm for obtaining the Bayes estimator, which is equivalent to setting a threshold and only letting those atoms whose realized atomic false discovery rate is below it to enter the Bayes estimator. We also provide alternate loss functions: Thus, we may either introduce linear constraints, which are equivalent to a regularization penalty, or consider the fractions of missed discoveries and false discoveries.

References

- Anderson, J. and V. Blair (1982): “Penalized maximum likelihood estimation in logistic regression and discrimination,” *Biometrika*, 69, 123–136.
- Baldi, P. and A. Long (2001): “A Bayesian framework for the analysis of microarray expression data: regularized t-test and statistical inferences of gene changes,” *Bioinformatics*, 17, 509–519.

- Benjamini, Y. and Y. Hochberg (1995): “Controlling the false discovery rate: a practical and powerful approach to multiple testing,” *Journal of the Royal Statistical Society B*, 57, 289–300.
- Bouton, C. and J. Pevsner (2002): “DRAGON View: information visualization for annotated microarray data,” *Bioinformatics*, 18, 323.
- Cui, X., J. Hwang, J. Qiu, N. Blades, and G. Churchill (2005): “Improved statistical tests for differential gene expression by shrinking variance components estimates,” *Biostatistics*, 6, 59–71.
- Do, K., P. Müller, and F. Tang (2005): “A Bayesian mixture model for differential gene expression,” *Applied Statistics*, 627–644.
- Efron, B. and R. Tibshirani (2002): “Empirical bayes methods and false discovery rates for microarrays,” *Genetic Epidemiology*, 23, 70–86.
- Garey, M. and D. Johnson (1979): *Computers and intractability: a guide to NP-completeness*, WH Freeman and Company, San Francisco.
- Gill, P., W. Murray, and M. Wright (1981): *Practical optimization*, London: Academic Press.
- Goeman, J. and P. Buhlmann (2007): “Analyzing gene expression data in terms of gene sets: methodological issues,” *Bioinformatics*, 23, 980.
- Gottardo, R., J. Pannucci, C. Kuske, and T. Brettin (2003): “Statistical analysis of microarray data: a Bayesian approach,” *Biostatistics*, 4, 597.
- Green, P. and B. Silverman (1994): *Nonparametric regression and generalized linear models: A roughness penalty approach*, New York: Chapman and Hall.
- Irizarry, R., C. Wang, Y. Zhou, and T. Speed (2009): “Gene set enrichment analysis made simple,” *Johns Hopkins University, Dept. of Biostatistics Working Papers*, 185.
- Lönnstedt, I. and T. Speed (2002): “Replicated microarray data,” *Statistica Sinica*, 12, 31–46.
- Mirnics, K., F. Middleton, A. Marquez, D. Lewis, and P. Levitt (2000): “Molecular characterization of schizophrenia viewed by microarray analysis of gene expression in prefrontal cortex,” *Neuron*, 28, 53–67.
- Mootha, V., C. Lindgren, K. Eriksson, A. Subramanian, S. Sihag, J. Lehar, P. Puigserver, E. Carlsson, M. Ridderstrale, E. Laurila, et al. (2003): “PGC-1 α -responsive genes involved in oxidative phosphorylation are coordinately down-regulated in human diabetes,” *Nat. Genet.*, 34, 267–273.
- Müller, P., G. Parmigiani, and K. Rice (2007): “FDR and Bayesian multiple comparisons rules,” *Bayesian statistics*, 8.
- Müller, P., G. Parmigiani, C. Robert, and J. Rousseau (2004): “Optimal sample size for multiple testing: the case of gene expression microarrays,” *Journal of the American Statistical Association*, 99, 990–1002.
- Newton, M. and C. Kendziorski (2003): “Parametric empirical bayes methods for

- microarrays,” in R. I. G. Parmigiani, E.S. Garrett and S. Zeger, eds., *The analysis of gene expression data: methods and software*, New York: Springer Verlag.
- Newton, M., C. Kendzioriski, C. Richmond, F. Blattner, and K. Tsui (2001): “On differential variability of expression ratios: improving statistical inference about gene expression changes from microarray data,” *Journal of Computational Biology*, 8, 37–52.
- Newton, M., A. Noueir, D. Sarkar, and P. Ahlquist (2004): “Detecting differential gene expression with a semiparametric hierarchical mixture method,” *Biostatistics*, 5, 155.
- Parmigiani, G., E. Garrett, R. Anbazhagan, and E. Gabrielson (2002): “A statistical framework for expression-based molecular classification in cancer,” *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, 64, 717–736.
- Parsons, D., S. Jones, X. Zhang, J. Lin, R. Leary, P. Angenendt, P. Mankoo, H. Carter, I. Siu, et al. (2008): “An Integrated Genomic Analysis of Human Glioblastoma Multiforme,” *Science*, 321, 1807.
- Smyth, G. (2004): “Linear Models and Empirical Bayes Methods for Assessing Differential Expression in Microarray Experiments,” *Statistical Applications in Genetics and Molecular Biology*, 3, 1027.
- Sotiriou, C., P. Wirapati, S. Loi, A. Harris, S. Fox, J. Smeds, H. Nordgren, P. Farmer, V. Praz, B. Haibe-Kains, et al. (2006): “Gene expression profiling in breast cancer: understanding the molecular basis of histologic grade to improve prognosis,” *JNCI Journal of the National Cancer Institute*, 98, 262.
- Storey, J. (2002): “A direct approach to false discovery rates,” *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 64, 479–498.
- Storey, J. (2003): “The positive false discovery rate: A Bayesian interpretation and the q-value,” *The Annals of Statistics*, 31, 2013–2035.
- Storey, J., J. Akey, and L. Kruglyak (2005): “Multiple locus linkage analysis of genomewide expression in yeast,” *PLoS Biology*, 3.
- Storey, J. and R. Tibshirani (2003a): “Sam thresholding and false discovery rates for detecting differential gene expression in dna microarrays,” in R. I. G. Parmigiani, E.S. Garrett and S. Zeger, eds., *The analysis of gene expression data: methods and software*, New York: Springer Verlag.
- Storey, J. and R. Tibshirani (2003b): “Statistical significance for genomewide studies,” *Proceedings of the National Academy of Sciences of the United States of America*, 100, 9440.
- Subramanian, A., P. Tamayo, V. Mootha, S. Mukherjee, B. Ebert, M. Gillette, A. Paulovich, S. Pomeroy, T. Golub, E. Lander, et al. (2005): “Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles,” *Proceedings of the National Academy of Sciences*, 102, 15545–15550.

Tavazoie, S., J. Hughes, M. Campbell, R. Cho, and G. Church (1999): "Systematic determination of genetic network architecture," *Nature Genetics*, 22, 281–285.

Wright, S. and J. Nocedal (2006): *Numerical optimization*, Springer.