

# Report on Independent Hypothesis Weighting

Xihan Qian

2024-03-01

# 1. Introduction

Numerous techniques have been devised for the analysis of high-throughput data to accurately quantify biological features, including genes and proteins. To ensure the reliability of discoveries, the false discovery rate (FDR) has become the dominant approach for setting thresholds. The methods for controlling FDR primarily rely on  $p$ -values, among which the Benjamini-Hochberg (BH) procedure (Benjamini and Hochberg 1995) and Storey’s  $q$ -value (Storey 2002) are popular.

However, Ignatiadis, Klaus, Zaugg, and Huber suggest that FDR methods based solely on  $p$ -values exhibit suboptimal power when the individual tests vary in their statistical properties. (Ignatiadis et al. 2016). When these methods focus exclusively on  $p$ -values, they overlook potentially relevant covariates. For example, in RNA-seq differential expression analysis, one such covariate could be the normalized mean counts of genes. Intuitively, genes with higher counts are likely to have greater power in detection compared to those with lower counts, and it would be optimal to include this relationship in the analysis.

To address this limitation, a novel approach known as independent hypothesis weighting (IHW) has been introduced. IHW improves the power of multiple hypothesis testing while controlling the FDR rate. The key innovation of IHW is that it recognizes that not all tests have the same power or the same prior probability of being true. It allows for the assignment of different weights to different hypotheses based on covariates that are predictive of the test’s power or its probability of being a true discovery. These covariates can be anything from the biological characteristics of the genes being tested to the technical aspects of the measurements. More explorations will be done in the following sections on this method.

## 2. Theories

There have been existing attempts to increase power by using covariates and this section starts off by introducing some known methods.

### 2.1 Weighted and Group weighted BH procedure

The weighted BH method has  $m$  hypotheses  $H_1, H_2, \dots, H_m$  and  $m$  weights  $w_1, w_2, \dots, w_m \geq 0$  that satisfy  $\frac{1}{m} \sum_{i=1}^m w_i = 1$ . After the weights are obtained, the BH procedure is applied on the modified  $p$ -values:  $\frac{p_i}{w_i}$ . In this scenario, selecting the weights before observing the  $p$ -values is essential, relying on prior knowledge or information indicating the likelihood of some hypotheses being true over others. This requirement presents a significant challenge, one that IHW seeks to address, and it utilizes the grouped weighted BH procedure (GBH). In this method, there are  $G$  groups where each  $X_i$  takes the same value within each group. GBH first estimates the proportion of null hypotheses by  $\hat{\pi}_0(g)$ , then weights the hypotheses proportionally to  $\frac{1 - \hat{\pi}_0(g)}{\hat{\pi}_0(g)}$ , and finally apply the BH procedure. However, the asymptotic theories in this method doesn’t function well when the number of hypotheses  $\frac{m}{G}$  is finite (Ignatiadis and Huber 2021). To resolve this issue, cross-weighting is used, where the idea is analogous to cross-fitting in regression settings, and this gives rise to the naive version of IHW.

## 2.2 Naive IHW and two-groups model

This version, also abbreviated as IHW-GBH since it is based off GBH, first divides the hypothesis tests into  $G$  groups based on the values of covariate  $X = (X_1, X_2, \dots, X_m)$ , with  $m_g$  number of hypotheses in the  $g$ -th group. It is also assumed that we have access to the  $p$ -values  $P = (P_1, P_2, \dots, P_m)$ , which is independent of  $X$  under the null. With this setup,  $\sum_{g=1}^G m_g = m$ . Then the weighted BH procedure is applied with each possible weight vector  $w = (w_1, w_2, \dots, w_G)$ , while the optimal  $w^*$  is the vector that lead to the most rejections. This method extends the BH procedure, making it pertinent to discuss the associated maximization problem. This problem stems from the two groups model (Efron 2008), which is a Bayesian framework that explains the BH procedure. Formally, assume that  $H_i$  takes values 0 or 1, and  $\pi_0 = \mathbb{P}(H_i = 0)$ . The distributions are as follows:

$$\begin{aligned} H_i &\sim \text{Bernoulli}(1 - \pi_0) \\ P_i \mid H_i = 0 &\sim U[0, 1] \\ P_i \mid H_i = 1 &\sim F_1 \end{aligned}$$

The marginal distribution for  $p$ -value  $P_i$  is then

$$P_i \sim F(t) = \pi_0 t + (1 - \pi_0) F_1(t)$$

With this, the Bayesian FDR becomes:

$$\text{Fdr}(t) = \mathbb{P}[H_i = 0 \mid P_i \leq t] = \frac{\pi_0 t}{F(t)}$$

A natural empirical estimator for the CDF would be the ECDF, and it can be written in terms of  $R(t)$ , which denotes the total number of rejections:

$$R(t) = m\widehat{F}(t) = \sum_{i=1}^m \mathbf{1}_{\{P_i \leq t\}}$$

Hence if  $\widehat{\pi}_0$  is an estimator of  $\pi_0$ ,

$$\widehat{\text{Fdr}}(t) = \frac{\widehat{\pi}_0 t}{\widehat{F}(t)} = \frac{\widehat{\pi}_0 m t}{R(t)}$$

If a conservative estimate is made:  $\widehat{\pi}_0 = 1$ , then

$$\widehat{\text{Fdr}}(t) = \frac{m t}{R(t)} \tag{1}$$

With this, the optimization problem is:

$$\text{maximize } R(t), \text{ s.t. } \widehat{\text{Fdr}}(t) \leq \alpha, t \in [0, 1]$$

The corresponding estimator in IHW-GBH is of the form

$$\widehat{\text{Fdr}}(t, \mathbf{w}) = \frac{mt}{R(t, \mathbf{w})} = \frac{\sum_{g=1}^G m_g w_g t}{R(t, \mathbf{w})}$$

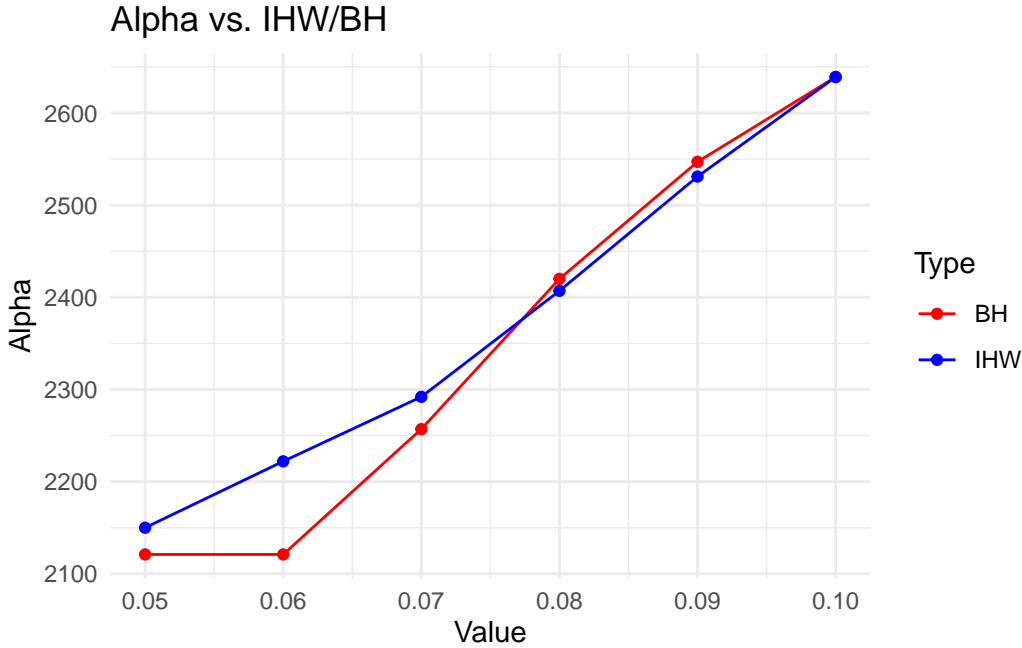
where  $R(t, \mathbf{w}) = \sum_{i=1}^m \mathbf{1}_{\{P_i \leq w_g t\}}$  is the number of rejections in bin  $g$ . Now the optimization problem is:

$$\text{maximize } R(t, \mathbf{w}), \text{ s.t. } \widehat{\text{Fdr}}(t, \mathbf{w}) \leq \alpha$$

However, with this approach, there are also some disadvantages including potential loss of Type I error, complications in solving the maximization problem, and its inability to scale when large number of tests are present. That is why modifications are made, leading to the IHW method.

### 2.3 IHW

In the first modification, the ECDF  $\hat{F}_g$  of the  $p$ -values in group  $g$  is replaced by the least concave majorant version called the Grenander estimator  $\tilde{F}_g$ . With this the maximization problem can be efficiently solved. The second modification involves randomly splitting the hypotheses into  $K$  folds, where  $K$  is usually taken to be 5. Then the maximization problem with the previous modification is applied to the remaining folds, leading to a weight  $\tilde{w} = (\tilde{w}_1, \dots, \tilde{w}_G)$ .



## References

- Benjamini, Yoav, and Yosef Hochberg. 1995. “Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing.” *Journal of the Royal Statistical Society: Series B (Methodological)* 57 (1): 289–300.
- Efron, Bradley. 2008. “Microarrays, Empirical Bayes and the Two-Groups Model.”
- Ignatiadis, Nikolaos, and Wolfgang Huber. 2021. “Covariate Powered Cross-Weighted Multiple Testing.” *Journal of the Royal Statistical Society Series B: Statistical Methodology* 83 (4): 720–51.
- Ignatiadis, Nikolaos, Bernd Klaus, Judith B Zaugg, and Wolfgang Huber. 2016. “Data-Driven Hypothesis Weighting Increases Detection Power in Genome-Scale Multiple Testing.” *Nature Methods* 13 (7): 577–80.
- Storey, John D. 2002. “A Direct Approach to False Discovery Rates.” *Journal of the Royal Statistical Society Series B: Statistical Methodology* 64 (3): 479–98.