OXFORD

## Gene expression

# FastLSU: a more practical approach for the Benjamini–Hochberg FDR controlling procedure for huge-scale testing problems

## Vered Madar[1] and Sandra Batista[2,]*

[1]Statistical and Applied Mathematical Sciences Institute, Research Triangle Park, NC 27709, USA and [2]Department of Computer Science, Princeton University, Princeton, NJ 08540, USA

*To whom correspondence should be addressed.
Associate Editor: Ziv Bar-Joseph

## Abstract

**Motivation:** We address a common problem in large-scale data analysis, and especially the field of genetics, the huge-scale testing problem, where millions to billions of hypotheses are tested together creating a computational challenge to control the inflation of the false discovery rate. As a solution we propose an alternative algorithm for the famous Linear Step Up procedure of Benjamini and Hochberg.
**Results:** Our algorithm requires linear time and does not require any *P*-value ordering. It permits separating huge-scale testing problems arbitrarily into computationally feasible sets or *chunks*. Results from the *chunks* are combined by our algorithm to produce the same results as the controlling procedure on the entire set of tests, thus controlling the global false discovery rate even when *P*-values are arbitrarily divided. The practical memory usage may also be determined arbitrarily by the size of available memory.
**Availability and implementation:** R code is provided in the supplementary material.
**Contact:** sbatista@cs.princeton.edu
**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

## 1 Introduction

In many fields the substantially increased scale of data available has resulted in a significant increase in the size of multiple hypotheses testing problems. In genetics, in particular, typical GWAS studies consist of $10^5 - 10^6$ SNPs (Hindorff *et al.*, 2015) while eQTL studies (Wright *et al.*, 2014; Xia *et al.*, 2012), newly advanced methylation studies (Smith *et al.*, 2014), and imaging studies (Stein *et al.*, 2010) usually start with $10^9$ tests. These testing problems are huge-scale as opposed to *large-scale* used by Efron (2004) to describe studies consisting of hundreds to thousands of hypotheses. It is preferable to control the false discovery proportion rather than the number of false positives for a huge-scale testing problem. Therefore the *FDR* or the *pFDR* approaches are favored and both tend to offer larger, more powerful sets of results than those yielded by the conservative *FWER* control. These huge-scale multiple hypotheses

testing problems create numerous computational challenges when many tests, say of the order $10^6$, are performed with all of the *P*-values of more or less equal importance. As a result some simpler testing procedures such as rigid *P*-value thresholds may be used that sacrifice power and correctness.

Alternatively tests may be separated or chunked into smaller sets or chunks that are more computationally feasible. Efron (2008) notes that the problem of separating hypotheses tests has not received great attention and warns of some pitfalls in chunking *P*-values, but focuses on grouping tests that share a biological property rather than arbitrary, computationally feasible chunks. Cai and Sun (2009) and later (Benjamini and Bogomolov, 2014) propose alternative solutions to Efron's grouping problem but do not address the problem of arbitrary, computationally feasible chunking. We confront the computationally feasible chunking problem for the

Benjamini–Hochberg false discovery rate (Benjamini and Hochberg, 1995). We show on data from Stranger's HapMap study (Stranger *et al.*, 2007; Xia *et al.*, 2012) that if results from separate tests are not combined correctly, there is considerable inflation of type I error, offer an explication for this occurrence, and propose our algorithm as a solution.

Consider a huge-scale testing problem of size $m$ where our goal is to select exactly $R \geq 0$ significant tests. Of the $R$ significant discoveries, exactly $V \geq 0$ tests will be false discoveries (i.e. truly non-significant tests that are declared significant). A common approach in multiple hypotheses testing problems is to control the family-wise error rate, $\mathrm{FWER} = \Pr(V \geq 1)$, the probability of selecting at least one false discovery. For huge-scale testing a more favorable alternative is to control the false discovery proportion, $\mathrm{FDP} = V/\max(R, 1)$, the proportion of truly false tests among the significant $R$. Some will prefer to control the positive FDR, $p\mathrm{FDR} = E(\mathrm{FDP}|R > 0)$, the expectation of the FDP when significant tests are selected, while others will opt to control the false discovery rate, FDR, the expectation of the FDP, $E(\mathrm{FDP})$. The FDR is always of a potentially smaller magnitude than the FWER and of the $p\mathrm{FDR}$ $(\mathrm{FDR} = p\mathrm{FDR} \cdot \Pr(R > 0))$. Yet, in reality for huge-scale testing, $\mathrm{FWER} \gg \mathrm{FDR}$ and sometimes, $p\mathrm{FDR} \approx \mathrm{FDR}$. Therefore both FDR and pFDR control approaches tend to offer larger, more powerful sets of results than those that might be offered by the conservative FWER control. For a further discussion about FWER, FDR and pFDR refer to Farcomeni (2004).

In huge-scale testing when the $m$ $P$-values are partitioned into chunks, it is challenging to control the FWER, pFDR or FDR over the entire collection of $m$ $P$-values. Controlling these error rates on a per chunk basis, if not done correctly, may interfere with the overall results by introducing more false discoveries. Although the same difficulty arises for the FDR and pFDR control (Benjamini and Bogomolov, 2014; Efron, 2008), it is easier to illustrate this for the FWER. Consider, for instance, an example of FWER control using the Bonferroni approach by collecting all $P$-values less than $\alpha/m$. Assuming that the number $m$ of $P$-values happen to be very large so that $m$ should be divided into $k$ chunks, each of of size $m_i$ so that $m = \sum_{i=1}^{k} m_i$. Applying Bonferroni in chunks of size $m_i$ will tend to select more significant results than applying it over the entire set of $m = \sum m_i$ $P$-values since $\alpha/m$ is less than $\alpha/m_i$. In the case of *FWER* control, using a fixed bound of $\alpha/m$ for all the chunks is theoretically preferred but often yields no significant results. A stricter constant cut-off on all sets of tests as suggested by Dudbridge and Gusnanto (2008) for GWAS was developed based on results from simulated GWAS. However such an ad-hoc approach eliminates from the entire multiple hypotheses testing problem any knowledge of the actual significance level $\alpha$ used.

### 1.1 The Benjamini–Hochberg linear step up procedure for controlling the false discovery rate

The Benjamini–Hochberg Linear Step Up (LSU) procedure is designed to control the FDR at desired signficance level $\alpha$ (Benjamini and Hochberg, 1995). As a result for huge-scale multiple hypotheses tests of equal importance, controlling a proportion of false discoveries, especially on the average, has increased power over procedures that control the FWER such as Bonferroni or a rigid cut-off bound such as $5 \times 10^{-8}$ suggested for GWAS (Dudbridge and Gusnanto, 2008). The larger the multiple hypotheses testing problem is the more powerful the LSU is over procedures that control the FWER.

While the LSU procedure is still one of the most cited procedures, its application had required sorting all $P$-values in decreasing

order to look for the largest $P$-value that satisfies a simple condition. In face of a huge-scale testing problem rather than apply LSU, some researchers had preferred to use harsh $P$-value cut-offs as mentioned above or to divide their huge-scale set of tests into computationally feasible smaller chunks and apply a multiple hypotheses testing procedure on each chunk selecting as the final significant results the union of the results in each of the chunks. Efron (2008) warns of the danger in such aggregation from the perspective of pFDR, pointing out that some chunks might have a larger proportion of significant results than others, and aggregating the significant results can yield misleading estimates. Moreover different chunking of tests might yield different sets of significant results. Analysis done by different groups of populations or chromosomes may not give the same number of significant tests as analysis that is applied on equal sized subsets; for example, it is well known that chromosome 6 (Mungall *et al.*, 2003) has a higher proportion of significant HLA SNPs than other chromosomes. We shall show in Section 2 that sorting the $P$-values, arbitrary thresholds, and arbitrary aggregation of results are not necessary and do not improve computational time efficiency compared to our algorithm.

### 1.2 A faster algorithm for LSU

Our alternative algorithm to the Benjamini–Hochberg LSU, **FastLSU**, performs linear scans instead of sorting $P$-values, but takes into account the overall size of the testing problem. FastLSU tiles the LSU procedure to give one global set of results that does not differ from applying LSU to the entire set of tests. Our algorithm addresses the same objective function as the original LSU. Our approach is provably faster than the conventional approach that relies on sorting $P$-values. It may also be used on arbitrary chunks of arbitrary size with an arbitrary space constraint in order to return the same set of significant results as those from applying LSU to the entire set of tests.

In the following section we address the difference between grouping and chunking tests and the difficulty in arbitrarily chunking tests by giving examples of inflation of type I error. In Section 3 we present FastLSU on a single set of tests and prove its equivalence to LSU, time efficiency and space efficiency. We present FastLSU on arbitrary chunks and also show its correctness and efficiency in Section 4. We offer suggestions for finalizing the report of significant tests in Section 5 and conclude with discussion in Section 6. R Code for an implementation of the algorithm is given in the supplementary materials.

## 2 An exercise of three HapMap groups and their chunking for stranger's *cis*-eQTL study

FastLSU controls the global FDR, not a family-based bound or a group-based bound. Our final results are not affected by the procedure of separating tests. This is an important distinction because chunking is arbitrary and based primarily on computational efficiency without any consideration for relationships between the hypotheses being tested. This is in contrast to group-testing procedures where for example, hypotheses may be grouped based on experimental knowledge such as all the tests from the same chromosome, the back and front parts of the brain (Efron, 2008), or different population groups in HapMap (Stranger *et al.*, 2007). In Section 3.2 we will show how the FastLSU algorithm can even improve the efficiency of a group-based controlling procedure. FastLSU is particularly suited to the current applications in genetics because we typically seek the significant set of SNPs or genes, and the family structure is usually of less importance than managing the

computational burden. Genetic family structures typically do not require any correction for family selection because each genetic family tends to contribute significant results of its own. (This is the case for the following example of 3 families of HapMap). For this reason in the remainder of this section, we will give motivating examples of applying FastLSU on the three groups of approximately $14 \times 10^6$ *cis*-eQTLs of Stranger's HapMap study (Stranger *et al.*, 2007; Xia *et al.*, 2012) in order to demonstrate the problems with arbitrary chunking without combining the results as FastLSU does.

Stranger *et al.* (2007) presents an eQTL study over 4 HapMap population samples: 30 Central Europeans(CEU), 45 Chinese (CHB), 45 Japanese(JPT) and 30 trios from Nigerians(YRI). To increase power each group is analyzed separately. We follow the recommendations of the SeeQTL website (Xia *et al.*, 2012) and consider the CHB and JPT together. We define *cis*-eQTLs as within 1 Mb upstream or downstream of a gene. Each population has a set of approximately 14 million *P*-values that were split into chunks of the following sizes: 1*M*, 900*K*, 800*K*, 700*K*, 600*K*, 500*K*, 250*K*, 100*K*, 50*K*, 25*K*, 10*K* and 5*K*. We contrast the differences in the number of significant results when the results are combined by taking the union of the results of LSU on each chunk versus those selected by FastLSU Algorithm 2 in Figure 1.

Applying the FastLSU yields 9228 significant results for CEU, 6497 for YRI and 33 507 for the CHB & JPT group. Most notable is that the results for FastLSU Algorithm 2 do not change across the chunk sizes whereas the alternative of taking the union of the results on each chunk increases not only the number of significant tests selected, but also the maximal *P*-value reported as the chunk size decreases. For example, applying FastLSU at level 10% for the complete chunk of approximately 14*M* *P*-values yields 9228 discoveries for the CEU. However, applying LSU on 1370 chunks of size 10*K* yields 16 925 significances which is equivalent to an overall FDR control of 23.75%. When 10*K* sized chunks are used for the CHB-JPT group, 41 587 *P*-values are selected as significant and this would require overall FDR control of $\alpha = 15.58\%$ For the YRI group when 10*K* sized chunks, there are 10 330 significant *P*-values and an $\alpha = 21.1\%$ would be required for this overall FDR control. This implies that performing the analysis using 1300–1400 chunks of size 10*K* rather than a single chunk inflates the type I error by 50%. While this is compelling experimental evidence, a more compelling explanation is that the
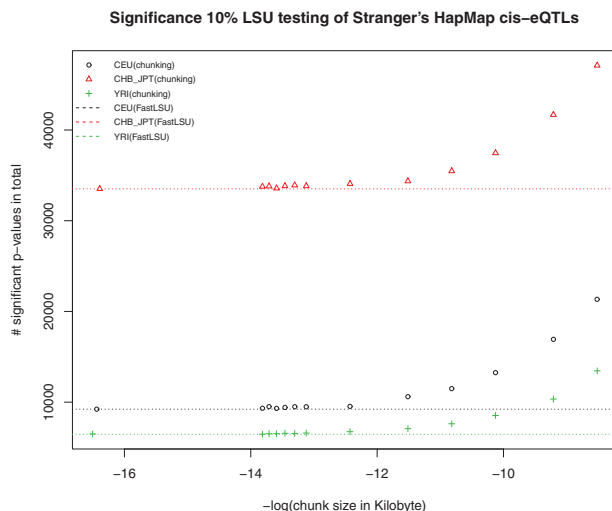
union of partial orderings on subsets of a set is in general not equal to the partial ordering on the entire set.

Moreover when the size of the chunks decreases, the significance interval is being divided into larger sub-intervals with each more likely containing more *P*-values spanning a greater range of values, so that the set of selected significant *P*-values will more likely contain a larger cut-off of *P*-values, and therefore a considerably larger value for the maximal significant *P*-value. The maximal value of significant *P*-values we obtain for the case of a single chunk is $6.7 \times 10^{-5}$ for CEU, $4.4 \times 10^{-5}$ for YRI and $2.5 \times 10^{-4}$ for CHB-JPT. For chunks of the size 100*K* (about 13 to 14 chunks), the max *P*-values are slightly higher especially for the CHB-JPT group: $1.7 \times 10^{-4}$ for CEU, $3.5 \times 10^{-4}$ for YRI and $1.2 \times 10^{-3}$ for CHB-JPT. For chunks of 10*K* size (about 130–140 chunks) the maximal *P*-values are considerably higher still: 0.0063 for CEU, 0.0034 for YRI and 0.011 for CHB-JPT. As we mentioned, there is no need to apply any correction for group or family selection as required by Benjamini and Bogomolov (2014) since all three groups in the HapMap example give significant results.

## 3 The FastLSU algorithm

The usual way to apply the LSU (Benjamini and Hochberg, 1995) at level $0 \leq \alpha \leq 1$ is to sort the *P*-values in descending order, $p_{(m)} \geq p_{(m-1)} \geq \cdots \geq p_{(1)}$. Starting from the largest *P*-value to the smallest, we need to look for the *k*th largest *P*-value that satisfies $p_{(k)} < k\alpha/m$. One may view the LSU algorithm (Benjamini and Yekutieli, 2001) as a search for optimal index $r$

$$r = \operatorname{argmax}\{i : p_{(i)} < i\alpha/m\}. \tag{1}$$

This observation motivates the following Fast Linear Step Up (FastLSU) algorithm that controls the FDR at level $\alpha$:

---

**Algorithm 1.** FastLSU for a single batch – without sorting the P-values

---

**Step 1.** *Start with $r_0 = m$ and count all P-values $< r_0\alpha/m$; let $r_1$ be their count.*

**Step 2.** *Do for $k = 2, \ldots, m$ Count all P-values $< r_{k-1}\alpha/m$; let $r_k$ be their count.*

**Step 3.** *Repeat step k for $k = k + 1$ until $r_k = r_{k+1}$ or $k = m$.*

---

**Example 1.** *We demonstrate Algorithm 1 on the example appearing in the original 1995 LSU paper (Benjamini and Hochberg, 1995) with $\alpha = 0.05$. Consider the 15 P-values:*

0.6528, 0.7590, 0.0298, 0.4262, 0.0459, 0.0278, 0.0001, 0.0019,

0.0004, 0.0201, 1.0000, 0.5719, 0.3240, 0.0095, 0.0344

*At the first step, we find 9 P-values $< 0.05$. That is, $r_1 = 9$. At the second step, $r_2 = 7$, 7 P-values are $< 9 \cdot 0.05/15 = 0.03$. At the third step, $r_3 = 5$ P-values are $< 7 \cdot 0.05/15 = 0.2333$. We stop with $r_4 = 4$ P-values that satisfy $< 4 \cdot 0.05/15 = 0.01333$. The selected P-values are $0.0001, 0.0019, 0.0004, 0.0095$, and the reader can check that these 4 P-values are exactly the ones selected using the original LSU P-value sorting algorithm.*

### 3.1 The equivalence to LSU and computational efficiency

THEOREM 1. *For a significance level $\alpha$, the Algorithm 1 maximizes the same objective function (1) that is used by the LSU. Therefore,*



**Significance 10% LSU testing of Stranger's HapMap cis–eQTLs**

**Fig. 1.** Varying chunk size over Stranger's HapMap populations. Dashed lines show the consistent result of applying the FastLSU as in Algorithm 2. The triangle, dot, and plus shapes represent results of applying the LSU under 10%

the FastLSU controls the FDR at level α and gives the same selected set of significant results as would be obtained by applying the Benjamini–Hochberg LSU FDR controlling procedure (Benjamini and Hochberg, 1995).

PROOF. Consider a batch of *P*-values which we will denote as $\mathcal{C}$. Let $0 \leq t \leq 1$ and define $\mathcal{S}(t : \mathcal{C}) = \{p < t : p \text{ is } P-\text{value} \in \mathcal{C}\}$ to be the number of *P*-values from $\mathcal{C}$ smaller than *t*. The set $\mathcal{S}\left(\frac{i\alpha}{m}\right)$ consists of all *P*-values from $\mathcal{C}$ that are smaller than $\frac{i\alpha}{m}$. Hence, for the case of a single batch, it is possible to verify that the search for

$$r = \text{argmax}\{i : i = \#\mathcal{S}(i\alpha/m)\} \qquad (2)$$

is equivalent to looking for the largest *r*th *P*-value satisfying $p_{(r)} \leq \frac{r\alpha}{m}$, as requested by (1).

THEOREM 2. The FastLSU Algorithm 1 requires O(m) time and O(m) space where m is the number of *P*-values to be considered.

The proof to Theorem 2 is given in the Appendix A with accompanying pseudocode for procedural language implementations. The main observation behind the linear time algorithm is that it is a constant time check where, in terms of what proportion of units of size α/m, a *P*-value is relative to α. It is then only a single scan to count the number of *P*-values that fall within each range relative to α and to find the range of *P*-values that satisfy the LSU condition. This is not commonly how statisticians consider LSU because we are not comparing the *P*-values directly to each other; we are comparing them by examining the range in which they fall relative to α and only implicitly to each other. We demonstrate the linear scans used in the linear algorithm in the following example.

**Example 2**. *To demonstrate this equivalence we apply FastLSU on the example appearing in the original 1995 FDR paper (Benjamini and Hochberg, 1995) with α = 0.05. Instead of sorting all 15 P-values we apply Algorithm 1 using 3 linear scans and find 4 significant P-values. Consider the 15 P-values from the example given in the original LSU paper (Benjamini and Hochberg, 1995) from Example 1. In the first linear scan, we label each P-value with what interval k, $k = \lceil \frac{mp}{\alpha} \rceil$, such that the P-value is at most kα/m respectively as follows:*

$$+, +, 9, +, 14, 9, 1, 1, 1, 7, +, +, +, 3, 11$$

*For P-values greater than α, a label + is used. During the same linear scan, we can also maintain counts for the number of P-values labeled with k from 1 to m = 15 as follows:*

$$3, 0, 1, 0, 0, 0, 1, 0, 2, 0, 1, 0, 0, 1, 0$$

*In the second linear scan, we check for k from m = 15 to 1 if the the number of significant P-values remaining, 9 in this case, minus the number of P-values in ranges examined so far equals to k as follows:*

$$9, 8, 8, 8, 7, 7, 5, 5, 4, 4, 4, 4$$

*This occurs when k = 4. On the third and final scan, we return as significant those P-values with interval, r ≤ 4. The selected P-values are* 0.0001, 0.0019, 0.0004, 0.0095 *as in Example 1.*

### 3.2 Applying FastLSU over families or groups of *P*-values of equal relevance

Efron (2008) distinguishes between the groups of imaging *P*-values that arise from the front brain and the back brain, and develops an empirical Bayes set-up to combine the significant *P*-values from the two families. Benjamini and Bogomolov (2014) extended Efron's idea of groups into the voxel families of MRI where each set of *P*-values from a specific voxel of a certain location are

treated as a separate, homogeneous group of relevance. For the first step each voxel group is analyzed by applying the LSU at level α. For the second step the number of significant groups is collected and LSU is performed again with another significance level $\alpha^*$ to correct for the selection of groups. If, for example, *S* groups out of *G* are shown to have at least one significant *P*-value for LSU of level α, $\alpha^*$ is set to $S\alpha/G$ for each of the groups (Benjamini and Bogomolov, 2014).

Given the equivalence between LSU and FastLSU, FastLSU may be used for each step in the approach of Benjamini and Bogomolov. The second step needs only be done for groups that have at least one significant result. Since only significant *P*-values for each group need to be considered, FastLSU may be applied to each group by using $\alpha^{**} = \frac{S r_g(\alpha)}{G m_g}\alpha$ where each group $g$ $(g = 1, 2, \ldots, G)$ contains $m_g$ *P*-values of which $r_g(\alpha)$ were declared significant of level α. Since FastLSU achieves efficiency by the simple recalling and correction of the constants *m* (or $m_g$) for the number of tests, a similar approach also works for Storey-Efron's positive FDR approach (Efron, 2004; Storey, 2002). To be more precise, the complete grid of λ values should be collected in each step and then the q-values should be corrected accordingly.

## 4 Applying LSU on arbitrary chunks of *P*-values

Algorithm 1 may be extended for when *P*-values are divided into an arbitrary number of chunks of arbitrary size. The size of the chunks may be determined by practical memory constraints or any other criteria. The iterative steps are applied on each chunk. The significant results from each chunk are combined using another iterative step to form the final set of significant results.

---

**Algorithm 2**. FastLSU for the case of two or more chunks of P-values

---

Given a set of m tests and P-values divided into n separate chunks of sizes $m_i$ P-values for $i = 1, \ldots, n$ such that $m = \sum_{i=1}^{n} m_i$:
- **Step 1**. On each of the chunks, count the number of P-values less than α. Denote this count as $r_i^{(1)}$.
- **Step $k + 1$**. On each of the chunks (starting from $k = 1$) count the number of P-values less than $\left(\sum_{i=1}^{n} r_i^{(k)}\right) \alpha /m$, and let $r_i^{(k+1)}$ be the count for the chunk at this step.
- **Repeat Step $k + 1$** until $k + 1 = m$ or $\sum_{i=1}^{n} r_i^{(k+1)} = \sum_{i=1}^{n} r_i^{(k)}$. Mark these P-values as significant under LSU.

---

An example of Algorithm 2 applied to the example appearing in the original 1995 FDR paper (Benjamini and Hochberg, 1995) is given in the Appendix A.

There are many ways one can alter Algorithm 2. For example, *P*-values that do not satisfy the condition in the iterative step that preserves the LSU may be either flagged or dropped. Alternatively Algorithm 1 may be applied to each chunk independently starting with the initial value for $r_0 = m - m_i + r_i^{(1)}$ and still tiling by the total number of tests, *m*. The advantage of this is that this may be done in parallel if desired. The results of the chunks may then be combined into a single set, if space permits, and Algorithm 1 applied to this set to return the final set of significant results. If space does not permit the results of the chunks to be combined, they may be combined up to the space limit. Then the iterative step of Algorithm 2 above may be applied to the resultant chunks starting by checking for *P*-values less than $m'\alpha/m$ where $m'$ is the size of the union of the results.

THEOREM 3. *For a significance level* $\alpha$, *and a collection of* $c$ *chunks the Algorithm 2 gives the same selected set of significant results as would be obtained by applying the Benjamini–Hochberg* (Benjamini and Hochberg, 1995) *LSU FDR controlling procedure over the set of all P-values from the c chunks.*

PROOF. [Proof for Algorithm 2] Suppose we have exactly $c$ disjoint chunks of $P$-values $\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_c$, with sizes $m_i$ ($i = 1, 2, \ldots, c$) such that $\mathcal{C} = \cup_{i=1}^{c} \mathcal{C}_i$ and $|\mathcal{C}| = m = \sum_{i=1}^{c} m_i$ $P$-values. Let $r_{\mathcal{C}}$ be the number of selected significant $P$-values that satisfy the LSU objective (1):

$$r_{\mathcal{C}} = \operatorname{argmax}\left\{r : r = \#\mathcal{S}\left(\frac{r\alpha}{m} : \mathcal{C}\right)\right\}. \quad (3)$$

Algorithm 2 can be applied to each $c$ chunks and the last step of Algorithm 2 is to finalize the selection. Accordingly for chunk $\mathcal{C}_i$ we search for the largest $r_i^*$ that satisfy the objective:

$$r_i^* = \operatorname{argmax}\left\{r : r = \#\mathcal{S}\left(\frac{(r + m - m_i)\alpha}{m} : \mathcal{C}_i\right)\right\}. \quad (4)$$

We claim that

$$r_{\mathcal{C}} \leq r_i^* + m - m_i, \quad \text{for any } i = 1, \ldots, c \quad (5)$$

This implies that each of the $P$-values selected significant from the search for (3) within chunk $\mathcal{C}_i$ must be selected while searching for argmax $r_i^*$ as in (4) since $p \leq r_{\mathcal{C}} \cdot \alpha/m$ implies $p \leq (r_i^* + m - m_i)\alpha/m$.

It is possible to verify that $p_{(k:\mathcal{C}_i)}$ the $k$th $P$-value in chunk $\mathcal{C}_i$ cannot be larger than, $p_{(k+m-m_i)}$, the $(k + m - m_i)$ th $P$-value in the union $\mathcal{C}$:

$$p_{(k:\mathcal{C}_i)} \leq p_{(k+m-m_i)} \quad \text{for all } k = 1, \ldots, m_i. \quad (6)$$

In particular, by the condition (4), let $p_{(r_i^*)}$ be the largest $P$-value in $\mathcal{C}_i$ that satisfies $p_{(r_i^*:\mathcal{C}_i)} < (r_i^* + m - m_i)\alpha/m$. When $m_i > r_i^*$, it follows that, for $k = 1, \ldots, m_i - r_i^*$,

$$(r_i^* + k + m - m_i)\alpha/m < p_{(r_i^*+k:\mathcal{C}_i)} \leq p_{(r_i^*+k+m-m_i)}. \quad (7)$$

(When $m_i = r_i^*$ (5) holds trivially since at most the entire set of tests may be selected as significant). To prove (5) by contradiction, let us compare between $p_{(r_{\mathcal{C}})}$ to $p_{(r_i^*+k+m-m_i)}$. If we assume that $r_{\mathcal{C}} > r_i^* + m - m_i$, then we can define a positive integer $k = r_{\mathcal{C}} - (r_i^* + m - m_i)$ for which following the inequality in (7) holds,

$$
\begin{aligned}
r_{\mathcal{C}} \cdot \alpha/m &= (r_i^* + k + m - m_i)\alpha/m \\
&< p_{(r_i^*+k:\mathcal{C}_i)} \\
&\leq p_{(r_i^*+k+m-m_i)} = p_{(r_{\mathcal{C}})}.
\end{aligned}
$$

However, this is in contradiction to the condition (3) that provides $p_{(r_{\mathcal{C}})} < r_{\mathcal{C}} \cdot \alpha/m$.

In conclusion we show that applying FastLSU over the global union of $P$-values give the exact same selection of significant $P$-values. We search for $r_{\mathcal{C}}^*$ over the union of the results on each chunk of size $\sum_{i=1}^{c} r_i^*$

$$r_{\mathcal{C}}^* = \#\mathcal{S}\left(\frac{r_{\mathcal{C}}^* \cdot \alpha}{m} : \mathcal{S}(r_1^* : \mathcal{C}_1) \cup \cdots \cup \mathcal{S}(r_c^* : \mathcal{C}_c)\right).$$

Since the inequality (5) ensures that $P$-values that were not selected under the condition (4) would have not been selected under the condition (3), we conclude that $r_{\mathcal{C}}^* = r_{\mathcal{C}}$. ¶

THEOREM 4. *The FastLSU Algorithm 2 requires* O(m) *time and* O(m) *space where m is the number of P-values to be considered. The practical memory usage may be restricted to an arbitrary limit for the largest chunk size,* $m^*$.

Proof for the running time and space limitations are shown in the Appendix A.

## 5 Useful tips for finalizing the report of significant *P*-values

If we assume that the $R$ tests were declared significant by applying either the LSU or FastLSU, we offer tips for finalizing the set of significant $P$-values. In the remainder of this section we explain how to protect the FDR control of FastLSU against dependency structures and when such correction is actually needed. We will also explain how to compute q-values (adjusted $P$-values) for the final results without keeping the entire set of $P$-values and show how to add a set of simultaneous confidence intervals for the significant test statistics while accounting for the selection effect.

### 5.1 Correcting against general case of dependence

The LSU procedure is conservative under the general type of positive regression dependence on subsets (PRDS) (Benjamini and Yekutieli, 2001), so applying the LSU at significance level of $\alpha$ always ensures FDR $\leq \alpha$ for PRDS $P$-values. The PRDS class contains a larger set of structures, among them, the independent case and any positively associated $P$-values such as $P$-values obtained from a two-sided $t$ test. Since the LSU and FastLSU are equivalent, applying the FastLSU at level $\alpha$ on PRDS $P$-values will control the FDR at level-$\alpha$, as well.

For other types of non PRDS dependency, such as in the case of pairwise comparisons, it is recommended to use the Benjamini–Yekutieli procedure (Benjamini and Yekutieli, 2001) that applies LSU using $\alpha^* = \alpha/c$ instead of $\alpha$ where $c = \sum_{k=1}^{m} 1/k \approx \ln m + 0.5772$. By the same argument the FastLSU under non PRDS will have FDR $\leq \alpha$ when applied under $\alpha^* = \frac{\alpha}{\ln m + 0.5772}$. As a matter of practice, we suggest to first perform FastLSU using the significance level $\alpha$. Then, if the $P$-values are non PRDS, correct the $R$ selected $P$-values by applying FastLSU again over the single batch consisting of $R$ $P$-values using $\alpha^{**} = \frac{R\alpha^*}{m}$. This approach provides FDR $\leq \alpha$ since $\alpha^* < \alpha$.

### 5.2 How to compute *q*-values to a selected subset of significant tests

When it is preferable to report $q$-values or adjusted $P$-values, we suggest how this may be done more efficiently. If we assume that $R$ tests were selected as significant, let the $P$-value $p_{(R)}$ be the largest $P$-value satisfying $p_{(R)}m/R < \alpha$. All $P$-values larger than $p_{(R)}$ were not selected as significant since $p_{(i)}m/i > \alpha$ for $i > R$ and have $q$-values $> \alpha$. Therefore it is sufficient to consider only the set of selected $R$ $P$-values. The LSU $q$-value, $q_{(i)}$, for the $P$-value, $p_{(i)}$ has the form (Yekutieli and Benjamini 1999)

$$
\begin{aligned}
q_{(i)} &= \min_{j=R, R-1 \ldots, i+1}\{p_{(i)}m/i, q_{(j)}\}, \quad \text{for } i < R \\
q_{(R)} &= p_{(R)}m/R, \quad \text{for } i = R.
\end{aligned}
\quad (8)
$$

From this we can see that the algorithms presented by (Yekutieli *et al.*, 1999) and (Storey, 2002) are $O(R\log R)$. One needs only sort the $R$ selected $P$-values in descending order and then beginning from largest $P$-value assign the corresponding $q$-value in a final linear scan recording the minimum $q$-value assigned thus far. The q-values will also be descending assigned in this way and there is no need to compare previous values except the minimum $q$-value thus far. The value for the $q$-value for $q_{(i)}$ will only change when it is less than the minimum seen thus far and the minimum $q$-values will span a range

of *P*-values until either a sufficient number of *P*-values have been covered or a sufficient range in *P*-values have been covered. Alternatively, one can use existing procedures (such as the function p.adjust in R software or PROC MULTTEST in SAS) to compute adjusted *P*-values for the set of *R* selected *P*-values and multiply the results by *R*/*m* to correct them.

## 5.3 Confidence intervals for selected subset of significant results

A less common approach in genetic studies is to report the significant test results by constructing a set of confidence intervals for the test statistics. While *P*-value is merely a measure of the magnitude of the test statistic, a confidence interval may offer the additional information about the dispersion of that magnitude. The selection adjusted confidence intervals (Benjamini and Yekutieli, 2005) offer an appropriate construction that corrects against the false coverage effect of selection. For a useful example see (Jung *et al.*, 2011) for how this method is used for the significant log-fold changes of RNA Microarrays.

## 6 Discussion

We presented an efficient algorithm to apply correctly the Benjamini–Hochberg Linear Step Up *FDR* controlling procedure in a huge-scale testing problem. Since we have shown that our algorithm requires only linear time, our algorithm is provably not any more computationally burdensome than even using a rigid Bonferroni cut-off for control. However, unlike the rigid Bonferroni cut-off, our approach ensures the FDR control at level α and this is a more powerful alternative to controlling the FWER, especially when the multiple testing problem is of huge-scale. Our approach is also scalable to any subsetting or chunking of the overall subset of *P*-values.

In addition we offered tips for performing the LSU over a huge-scale multiple hypotheses testing problem, such as showing how to correct for dependency or how to compute q-values directly from the subset of LSU significant results. We hope that these algorithms and tips offer better insight into the huge-scale testing problem rather than just black-box solutions. We encourage altering the steps in performing Algorithm 2, for instance, either by applying it sequentially or in parallel or a mixture of both.

The amount of inflated type I error we observed during the exercise on different chunk sizes of the HapMap populations strongly suggests the need for greater diligence in correctly separating hypotheses tests, whether the objective is to control the FWER, the FDR, or pFDR. This observation was also reported by Efron (2008), but a full investigation on real data with decreasing chunk sizes was not performed. Our suggested approach solves this for the case of the Benjamini–Hochberg FDR. As we have alluded, a similar approach can be adopted to control correctly the Efron–Storey's pFDR (Storey, 2002) and this remains open for further research. In addition the severe phenomena of inflated type I error we observed suggests more care may be required in reporting and interpreting results in genetics literature especially in the case of GWAS and eQTL studies. Namely to properly understand the significance results there is a need for consistent consideration of the algorithms or software used for controlling and separating the hypotheses tests and for recording the chunk sizes used for a study.

## References

Benjamini,Y. and Bogomolov,M. (2014) Selective inference on multiple families of hypotheses. *J. R. Stat. Soc. Ser. B*, **76**, 297–318.

Benjamini,Y. and Hochberg,Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B*, **57**, 289–300.

Benjamini,Y. and Yekutieli,D. (2001) The control of the false discovery rate in multiple testing under dependency. *Ann. Stat*, **29**, 1165–1188.

Benjamini,Y. and Yekutieli,D. (2005) False discovery rate-adjusted multiple confidence intervals for selected parameters. *J. Am. Stat. Assoc*, **100**, 71–93.

Cai,T.T. and Sun,W. (2009) Simultaneous testing of grouped hypotheses: finding needles in multiple haystacks. *J. Am. Stat. Assoc*, **104**, 1467–1481.

Dudbridge,F. and Gusnanto,A. (2008) P-value less than say $5 \times 10^{-8}$ can be regarded as convincingly significant. *Genet. Epidemiol.*, **32**, 227–234.

Efron,B. (2004) Large-scale simultaneous hypothesis testing: the choice of a null hypothesis. *J. Am. Stat. Assoc*, **99**, 96–104.

Efron,B. (2008) Simultaneous inference: when should hypothesis testing problems be combined? *Ann. Appl. Stat*, **2**, 197–223.

Farcomeni,A. (2004) *Multiple Testing Procedures Under Dependence, With Applications*, Ph.D. Thesis, Dipartimento di Statistica, Probabilita' e statistiche applicate, Universita' di Roma "La Sapienza".

Hindorff,L.A. *et al.* (2015) *A Catalog of Published Genome-Wide Association Studies*. Available at: www.genome.gov/gwastudies. (11 January 2015, date last accessed).

Jung,K. *et al.* (2011) Reporting FDR analogous confidence intervals for the log fold change of differentially expressed genes. *BMC Bioinf.*, **12**, 288.

Mungall,A.J. *et al.* (2003) The DNA sequence and analysis of human chromosome 6. *Nature.*, **425**, 805–811.

Smith,A.K. *et al.* (2014) Methylation quantitative trait loci (meQTLs) are consistently detected across ancestry, developmental stage, and tissue type. *BMC Genomics*, **15**, 148.

Stein,J.L. *et al.* (2010) Voxelwise genome-wide association study (vGWAS). *Neuroimage*, **53**, 1160–1174.

Storey,J. (2002) A direct approach to false discovery rates. *J. R. Stat. Soc. Ser. B*, **64**, 479–498.

Stranger,B.E. *et al.* (2007) Population genomics of human gene expression. *Nat. Genet.*, **39**, 1217–1224.

Wright,F.A. *et al.* (2014) Heritability and genomics of gene expression in peripheral blood. *Nat. Genet.*, **46**, 430–437.

Xia,K. *et al.* (2012) seeQTL: a searchable database for human eQTLs. *Bioinformatics*, **28**, 451–452.

Yekutieli,D and Benjamini,Y. (1999) Resampling-based false discovery rate controlling multiple test procedures for correlated test statistics. *J. Stat. Plan. Inference*, **82**, 171–196.

# A1. PROOFS AND CODES

## A1.1 Proof of theorem 2

PROOF. [Proof of Theorem 2] We present the algorithm here in three steps with running time for each step and accompanying pseudocode.

**Bin.** Classify all *P*-values into the bins of the interval, $[0, \alpha]$ each of size $1/m$ and keep a total of all *P*-values with value less than $\alpha$, $m^*$. Label each *P*-value with a value $k$ such that $k = \left\lceil \frac{mp^*}{\alpha} \right\rceil$ or $k = k + 1$ if $\frac{mp^*}{\alpha}$ equals $\left\lceil \frac{mp^*}{\alpha} \right\rceil$ where $\alpha$ is the significance level and $p^*$ is the given *P*-value being labeled. If $p^*$ is labeled with bin $k$ less than or equal to $m$, increment the count for bin $k$ and increment current *P*-value count, $m^*$. If a $p^*$ has label $k$ greater than $m$, it may be filtered, so no counts need to be incremented for such *P*-values (although they can be labeled with arbitrarily large values for $k$). *Running time and memory:* Labeling each *P*-value and incrementing the labeled bin count requires only constant time and a single pass through the *P*-values. In addition to storing the *P*-values, the *P*-value labels and bin counts also need to be stored, also requiring $O(m)$ space each, and a variable for the current *P*-value count, $m^*$.

**Accumulate.** In this step we find the $r^*$ the significant bin to return all *P*-values in bins less than or equal to this bin as significant. To do so, starting from the highest labeled bin's count, i.e. for $m$, keep a partial sum of the total number of *P*-values in the bins thus far. If the current bin's has a non-zero bin count and its value is equal to $m^*$ minus then the current partial sum, then return the current bin as $r^*$ the significant bin. *Running time and memory:* This step can be done in a single scan of the bin counts and only requires additional variables for the significant bin, $r^*$ and the partial sum.

**Return.** Return as significant all the *P*-values that were labeled with a $k$ less than or equal to $r^*$. *Running time and memory:* This requires only a single scan of the *P*-value labels and no additional space.

## A1.2 Pseudocode for proof of theorem 2

```
Let p_vals[m] be an array of P-values,
p_val_labels[m] be array of bin labels for P-values,
count_bins[min] be an array of bin counts,
m the total number of P-values, and alpha the
significance level.
Initialize p_val_labels and count_bins to 0.
sig_pvals = 0; #number of significant P-vals
for i = 1 to m {
  p* = p_vals[i];
  k = ceil(mp*/alpha); #Bin
  if (ceil(mp*/alpha) equals mp*/alpha) increment k;
  p_val_labels[i] = k;
  if (k less than or equal to m) {
    increment count_bins[k];
    increment sig_pvals;
  }
}
r* = 0; #significant bin
total_p_vals = 0; # partial sum of p_values
for i = m to 1 {
  add count_bins[i] to total_p_vals; #Accumulate
  if (sig_pvals - total_p_vals equals i) {
    r* = i and break;
```

```
  }
}
for i = 1 to m {
  if (p_val_labels[i] less than or equal to r*)
    #Return
    return p_vals[i] as signficant;
}
```

## A1.3 Example for Algorithm 2

**Example 3.** To demonstrate the Algorithm 2 consider the 15 *P*-values appearing in Example 1 and also in the original LSU paper (Benjamini and Hochberg, 1995). Further assume that for some reason the 15 *P*-values are divided into two chunks. The first chunk, say $\mathcal{C}_1$, consists of the first 8 *P*-values:

$$\mathcal{C}_1 = \{0.6528, 0.7590, 0.0298, 0.4262, 0.0459, 0.0278, 0.0001, 0.0019\},$$

*and the second chunk,* $\mathcal{C}_2$

$$\mathcal{C}_2 = \{0.0004, 0.0201, 1.0000, 0.5719, 0.3240, 0.0095, 0.0344\}.$$

Algorithm 2 applied for $\alpha = 0.05$ has three steps: In Step 1, we scan the 8 *P*-values in $\mathcal{C}_1$ and seek for the largest *P*-value, $p_{(i:\mathcal{C}_1)}$ satisfying $< (i + 15 - 8)\alpha/15$. This can be done in a similar manner to applying Algorithm 1. First scan gives 5 *P*-values $< (8 + 15 - 8)0.05/15 = 0.05$. Second scan gives 4 *P*-values $< (5 + 15 - 8)0.05/15 = 0.04$ and stops with the 4 *P*-values $< (4 + 15 - 8)0.05/15 = 0.0367$. The selected *P*-values for the first step are

$$\mathcal{S}(0.0367 : \mathcal{C}_1) = \{0.0298, 0.0278, 0.0001, 0.0019\},$$

In Step 2, we scan the 7 *P*-values in $\mathcal{C}_2$ and seek for the largest $p_{(i:\mathcal{C}_2)} < (i + 15 - 7)0.05/15$. First scan gives 4 *P*-values $< (7 + 15 - 7)0.05/15 = 0.05$ and immediately stops with these 4 *P*-values $< (4 + 15 - 7)0.05/15 = 0.04$. The selected *P*-values are

$$\mathcal{S}(0.04 : \mathcal{C}_2) = \{0.0004, 0.0201, 0.0095, 0.0344\}.$$

Next, we follow the last step of Algorithm 2 which is a combination step and applies on the collection of 8 *P*-values selected at the former steps:

$$\{0.0298, 0.0278, 0.0001, 0.0019, 0.0004, 0.0201, 0.0095, 0.0344\}.$$

All 8 selected *P*-values are clearly $< 0.05$. Second, 5 *P*-values $< 8 \cdot 0.05/15 = 0.0267$, and third scan finds out 4 *P*-values $< 5 \cdot 0.05/15 = 0.0167$ that also satisfy $< 4 \cdot 0.05/15$. The selected *P*-values are, again, 0.0001, 0.0019, 0.0004, 0.0095.

## A1.4 Proof of theorem 4

Proof. [Proof of Theorem 4] For $m$ *P*-values arbitrarily divided into $n$ chunks of size $m_c$ for $c = 1, ..., n$ such that the maximum chunk size is $m^*$, we show that the algorithm is still linear in $m$ and never uses more than $m^*$ space.

**Bin.** This step is as for Algorithm 1. It is important to note that the binning is done relative to $m$ and not the size of the chunk. The only important difference is that bin sums are not maintained because of the $m^*$ space limitation. Labels need not be stored either. This step also counts the sum, $m'$, of all *P*-values across all groups that are less than $\alpha$. *Running time and memory:* This requires a linear scan. A count variable can be kept for each group in order to get $m'$.

**Accumulate.** For each group, find the bin sums for the largest $m^*$ partition of bins not covered yet, i.e. find bin sums for bins $m' - j$

$m^* + 1$ to $m - (j-1)m^*$ for $j = 1, ..., n$ and we do both of the following before incrementing $j$. Accumulate bin sums across the chunks. This can be done in a linear scan of the chunks and a single array accumulating bin sums across chunks. Finding $r*$ is as in Algorithm 1 Step 2. However, now only $m^*$ bins may be checked at a pass before needing to increment $j$. The subtotal of $P$-values

counted thus far is maintained after $j$ is incremented. *Running time and memory*: This step must be repeated at most $n$ times and requires a linear scan of the data. At any point at most $m^*$ space plus several count variables are used.

**Return.** This step is as in Algorithm 1.