

1 A Practitioner's Guide to Multiple Testing
2 Error Rates

3 Jonathan Rosenblatt

Department of Statistics and Operations Research,

The Sackler Faculty of Exact Sciences,

Tel Aviv University

Israel

4 June 25, 2013

5 **1 Introduction**

6 It is quite common in modern research for a researcher to test many hypothe-
7 ses. The statistical (frequentist) hypothesis testing framework does not scale
8 with the number of hypotheses in the sense that naïvely performing many
9 hypothesis tests will probably yield many false findings; “false” in the sense
10 they will not be replicated. Indeed, classical statistical “significance” is evi-
11 dence for the presence of a signal within the noise expected in a single test,

12 not in a multitude where the noise levels are higher. Strong evidence of signal
13 assuming one noise level, can easily be considered as no evidence of signal
14 under a higher noise level. For protection from an uncontrolled number of
15 erroneous findings, a researcher has to consider the type of errors, or non-
16 replications, he wishes to avoid. The researcher can then select the adequate
17 procedure for that particular error type and data structure, or alternatively
18 estimate that error type for a particular set of candidate findings.

19 In practice, the selection of the proper error rate might cause the re-
20 searcher some confusion. This point was made at the 2009 Multiple Com-
21 parisons conference in Tokyo [2, Section 4.4], demonstrated in the following
22 question from the statistics Questions & Answers web site *Cross Validated*¹ :

23 *I am testing many (500,000) genetic variants, and the tests*
24 *are FDR corrected and give me a q-value. Normally I would just*
25 *call everything with $q < .05$ significant. But in this case I am test-*
26 *ing those same genetic variants in two other related experiments*
27 *(not using exactly the same individuals, but the samples may over-*
28 *lap). What to do? Would changing the significance threshold for*
29 *q to $.05/3 = .0167$ be an option?*

30 This particular example is further discussed in Section 3.6.

31 To offer guidance, we review possible error types for multiple testing
32 (sec 2) and demonstrate them with some practical examples (sec 3) which

¹ See <http://stats.stackexchange.com/questions/26588/multiple-fdr-corrected-experiments-using-the-same-data>. Accessed on Apr 20, 2013

clarify the formalism of sec 2. Finally, in appendix A, we include some notes on the software implementations of the methods discussed.

A multiplicity control procedure (e.g. Bonferroni, Benjamini-Hochberg, ...) is a data manipulation process— an algorithm— that guarantees that a preselected error rate is no larger than a preselected value. A typical procedures will actually offer guarantees vis-à-vis several error measures simultaneously. The emphasis of this manuscript is however on the error rates, and not on the multiplicity control procedures themselves.

For the purpose of selecting the appropriate procedure consult your favorite software’s documentation (see our appendix A). Alternatively, Farcomeni [15] or more recently Goeman and Solari [19], can serve as references. As the focus of this paper is the error measures, p-value adjustment, simultaneous confidence intervals, and error estimation will not be discussed. The reader is referred again to [15] or [19] as possible references.

2 Measures of Error

2.1 Family Wise Error Rates

Consider the testing of several null hypotheses against their respective research (alternative) hypotheses. The Family Wise Error Rate (FWER) is the frequency of experiments in which a false rejection of some null hypothesis will occur; put differently, the probability of a false finding.

As is customary in single hypothesis testing, a FWER level of $\alpha = 0.05$ is

54 often used and sometimes even required, as in drug registering experiments.

Table 1 introduces the nomenclature which has become standard in the multiple comparisons community and will be referenced throughout this article. Following this notation, the FWER is defined as

$$Prob(V \geq 1)$$

55 where $Prob(.)$ denotes the relative frequency over repeated experiments.

	Claimed nonsignificant	Claimed Significant	Total
Null	U	V	m_0
Nonnull	T	S	m_1
Total	$m - R$	R	m

Table 1: Classification of types of decisions made

56 2.1.1 Weak and Strong Control the Family Wise Error Rate

57 The probability of any particular inference procedure of making a false finding
58 depends on the existence of true effects. FWER control in the “weak” sense
59 refers to procedures which guarantee controllable FWER when there are no
60 true effects at all, i.e., when all null hypotheses are correct.

61 Detecting the existence of *any* phenomena ($m_1 > 1$) is a simpler task that
62 actually identifying these phenomena. Lack of weak FWER control means
63 that we have no error guarantees even regarding this simple task. As men-
64 tioned in Section 1, a multiplicity control procedure might offer guarantees
65 with respect to several error measure simultaneously. Weak FWER control

66 should, and typically is, a minimal requirement.

67 FWER control in the “strong” sense is the complementing concept, re-
68 ferring to procedures which guarantee FWER control even in the presence
69 of some true effects, i.e., when not all null hypotheses are correct. As their
70 names imply, “strong” control is the stricter criterion, which entails “weak”
71 control.

72 2.2 False Discovery Rates

The False Discovery Rate (FDR), first introduced by Benjamini and Hochberg [5], is the ratio between false discoveries and total discoveries, averaged over replicated experiments. Denoting the (unknown) False Discovery Proportion in a particular experiment as $FDP = V/R$, and setting the convention that no discoveries signify no errors ($R = 0 \Rightarrow FDP = 0$), the FDR can be now defined as:

$$E(FDP)$$

73 where $E(.)$ denotes the average over all possible experimental results.

74 **Remark 2.1.** The FDR error rate has become synonymous with the Benjamini-
75 Hochberg procedure presented in [5] . This is plain wrong and confusing.
76 Benjamini-Hochberg is a procedure that does indeed offer FDR control in
77 particular setups, but it is only one of many.

78 **Remark 2.2.** In the context of FDR, there need not be an $\alpha = 0.05$ con-
79 vention. Researchers are free to choose the level of error they see adequate

80 for their particular research. Obviously, an error level of $\alpha = 0.5$ might
81 be hard to defend from critique. Having said that, since FDR control does
82 offer FWER control in the weak sense, the $\alpha = 0.05$ convention might be
83 justifiable.

84 2.3 Other Measures of Error

85 FWER and FDR are the most commonly used, but by no means the only
86 measures of error. Since many error measures are merely an average over
87 replications of the experiment, many other error measures can be considered
88 by replacing the error function to be averaged. Denoting by $E(C)$ the aver-
89 age, over all possible experimental outcomes, of some error C . We now see
90 that FWER and FDR simply set $C = I_{\{V \geq 1\}}$ and $C = FDP$ respectively.
91 Some other measures of error are:

- 92 • Per Family Error Rate (PFER): Where $C = V$.
93 This measure is the simple (expected) number of erroneous discoveries.
- 94 • Per Comparison Error Rate (PCER): Where $C = V/m$.
- 95 • k -FWER [29]: Where $C(k) = I_{\{V \geq k\}}$.
96 This measure is the relative frequency of the making of no less than k
97 erroneous discoveries.
- 98 • False Discovery Exceedance (FDX)[18]: Where $C(\gamma) = I_{\{V/R \geq \gamma\}}$.
99 This measure was motivated by the fact that FDR keeps the proportion

100 of false discoveries small, but only *on average*. In extreme scenarios,
 101 FDR does not exclude the possibility of making more than $FDP >$
 102 α mistakes in *almost all* experiments. FDX targets these scenarios
 103 explicitly, by allowing it to happen with a small probability, and is
 104 thus more conservative than FDR.

105 Other measures of error which are not simple averages over replications
 106 of the experiment include, but are not limited to:

- 107 • Positive FDR or pFDR [27] or FDR_{-1} [1] : Defined as $E(V/R; R > 0)$.
 108 This measure is essentially the proportion of false findings (within all
 109 findings), but averaged, not on all possible experiments outcomes, but
 110 only on those which actually return findings. It was motivated by the
 111 observation that the FDR of a procedure might be very low merely
 112 because in many events it returns no findings, even if it makes many
 113 mistakes when it does indeed return findings (see [27] for a example).
 114 On the other hand, if all the null hypotheses are true, thus all find-
 115 ings are false, one would want an error measure to coincide with the
 116 probability of a false finding (weak FWER). pFDR does not enjoy this
 117 attribute.

- 118 • Marginal FDR or mFDR [28] or Fdr [13] or FDR_{+1} [1]: Defined as
 119 $E(V)/E(R)$.

120 While not very interesting for itself, this error measure gained popu-
 121 larity since it is mathematically tractable and approximates the FDR

122 when many independent hypothesis are being tested.

123 We conclude by noting that FWER, FDR, pFDR and mFDR are (cur-
124 rently) by far the most popular. So much so that it is actually hard to find
125 published applied research using any other.

126 **2.4 Choosing Your Family**

127 All the previous error measures are defined for a family of hypotheses that is
128 known, and indirectly assumed this family is all the hypotheses being tested
129 in an experiment. This need not be the case, and defining the family of
130 relevance may be a non-obvious part of the researcher’s work. The examples
131 in Section 3 include some trivial scenarios, in the sense that the family of
132 hypotheses is clear. The section also includes some examples where the family
133 is not trivial (see 3.7) and its choice will depend on the scientific statement
134 in mind.

135 **3 Examples**

136 **3.1 Tukey’s Psychological Exams**

137 In his 1953 unpublished paper: “The Problem of Multiple Comparisons” [4]
138 and later, when lecturing at Princeton University [12], Prof. John Tukey
139 would tell a motivating tale about a young psychologist. After administering
140 250 tests he finds that 11 were significant at the 0.05 level. A null hypothesis

141 being that a given test does not differentiate between his groups of interest, a
142 (significant) rejection of this null means a test does actually differentiate the
143 groups. After the initial feeling of satisfaction he consults a senior researcher,
144 only to discover his findings are rather poor, since one would expect 12.5
145 significant tests due to chance alone. Having only 11 significant results is
146 actually disappointing.

147 With this new understanding, our psychologist now has to decide how
148 he can protect himself from false findings. Say the tests consist of new
149 candidate clinical diagnostics for condition X. Making an error means that a
150 test will be used to diagnose X while it actually cannot distinguish between
151 healthy and X. Since this is unacceptable for our psychologist, he will want
152 an inference procedure that controls the FWER in the strong sense. This
153 will also guarantee protection in the case that no test differentiates between
154 healthy and X, i.e., weak FWER control, which can actually be a question
155 of interest for itself.

156 Now consider a different scenario: The tests check for differences in per-
157 sonality attributes between genders. Making an error means that the psy-
158 chologist might believe male and female differ in a way they actually do not.
159 The researcher does not consider this a serious mistake, as long as many other
160 true differences are discovered. In this setup, the researcher should control
161 the FDR, FDX or pFDR. Allowing for some mistakes will allow the researcher
162 to enjoy a sensitivity gain compared to FWER-controlling-procedures. The
163 interpretation of the findings should be done in accord with the error measure

164 employed.

165 **3.2 ANOVA**

166 In their 1999 paper, Williams, Jones, and Tukey, analyze the National As-
167 sessment of Educational Progress (NAEP) 1990 and 1992 data. This data
168 consists of the average eighth grade mathematics proficiency scores for the
169 34 states that participated in both 1990 and 1992 NAEP Trial State Assess-
170 ment (TSA). Comparisons are made between regions (Central, Northeast,
171 Southeast and West), years (1990,1992), states nested within regions, and
172 the Year \times Region interaction. In this case a null hypothesis means there is
173 no difference in the proficiency score between sub-groups, and its rejection
174 meaning there is indeed such a difference.

175 We start by noting that this one study, provides four families of hypothe-
176 ses. Indeed, a falsely discovered difference between regions, as an example,
177 is of no concern when comparing years, states or regional changes.

178 We also note that in their paper, the authors actually compare between
179 procedures controlling the FWER and the FDR. They do not offer a justifi-
180 cation for the preference of any measure over the others, so we will offer one
181 of our own. We will remark however, that their bottom line is unorthodox
182 in the context of ANOVA:

183 *Each of the three authors believes that the B-H procedure is the*
184 *best available choice*

185 So why FWER? If, for example, the study is analyzed in the context of
186 discrimination– having no policy implications but rather a possible stigma-
187 tizing effect– the researcher might wish to refrain from any falsely discovered
188 differences between states, regions etc. If, on the other hand, intervention
189 policies are the context, then power is a major concern. Missing a difference
190 might mean policy makers are left unaware of the differences to be addressed.
191 This context requires a less stringent error criterion than FWER, leading to
192 the authors’ stated preferences.

193 **3.3 Functional Magnetic Resonance Imaging**

194 Consider now the case of the neuroscientist, trying to locate the brain regions
195 responsive to visual stimuli. He has scanned a dozen subjects or so in the
196 Magnetic Resonance Imaging (MRI) machine and recorded the brain’s acti-
197 vation² in response to the stimuli. To be precise, he measured the activation
198 level at *each* of several thousand brain locations, called Volumetric Picture
199 Elements (voxels); their exact number depending on the resolution of the
200 MRI scan . With the measured activation levels in hand, the researcher can
201 compute their correlation to the stimulus given. If the voxel-wise measure-
202 ment is correlated with the stimulus, the location is considered “active”. We
203 see that localizing activation actually consists of performing many local hy-
204 pothesis tests: the null hypothesis of no correlation to the stimulus is tested
205 at each voxel, and its rejection meaning a responsive location has been found.

² He actually measured the blood oxygenation level. Details can be found in [21].

206 Returning to multiplicity error rates; an error would mean the researcher
 207 declared a voxel as responsive when it actually is not. This does not seem
 208 like a terrible mistake to make, so the researcher should probably protect
 209 himself from large proportions of errors, and not from the making of one
 210 single error. FWER, k-FWER and Per Family Error Rate are thus excluded.
 211 Per Comparison Error Rate seems like a possible candidate, but it is very
 212 liberal. One can actually gain power by including many “junk” hypotheses,
 213 say, by including the air surrounding the head in the family. To see this,
 214 consider the case of infinitely many hypotheses tested. The proportion of
 215 errors will trivially be smaller than any α we pick.

216 Our researcher is thus left with FDR and FDX as candidates for measure-
 217 ment of error. In the case the researcher has no clear favorite from within
 218 these two measures, a possible consideration at this stage might be the avail-
 219 ability, simplicity and power of controlling procedures. These considerations
 220 give preference, at the time of writing, to FDR over FDX.

221 **Remark 3.1.** For fairness it should be stated that “cheating” with FDR and
 222 FDR is possible, by augmenting the family with some obviously *false* null
 223 hypotheses [16]. It is our view that “cheating” the FDR or FDX does re-
 224 quire more effort and malicious intent than analyzing a needlessly large brain
 225 volume. We thus do not qualify these two “cheats” as equally problematic.

226 3.3.1 Functional Magnetic Resonance Imaging– Cluster Level In- 227 ference

228 Return to the neuroscientist from sec 3.3. Recalling that the voxels are
229 arbitrary volume units, defined by the technology of the MRI and not by
230 entities of interest for inference, he decides that a more interesting entity is
231 a mass of contiguous activations. He thus decides that he is interested in
232 spatially contiguous regions with activation larger than “7” (in some scale).
233 These regions are known as “excursion regions”, “exceedance sets”, “blobs”
234 and possibly other names. After scanning a subject, he realizes there are
235 30 contiguous regions which exceed 7. Conscious that some are due merely
236 to chance variation, and knowing enough probability theory, he computes a
237 p-value for the observed volume (exceeding 7), in each of the 30 regions. If
238 he rather not make any mistakes, he can control the FWER of the regions.
239 Namely, controlling the probability of declaring any inactive regions as active.
240 This is indeed the approach implemented in several brain analysis software
241 packages, particularly SPM (<http://www.fil.ion.ucl.ac.uk/spm/>).

242 Alternatively, if the researcher wishes to allow for some slack and accept
243 false regions– as long as their proportion is not too high– he should use FDR
244 or FDX control. Alas, the FDR defined in Section 2.2 assumes an a priori
245 fixed and known number of hypotheses being tested (m). The number of
246 excursion regions is data dependent, thus random and a priori unknown.
247 Extensions of the FDR for the random-number-of-hypothesis case do exist.
248 A rigorous exposition can be found in Siegmund et al. [24]. Note however,

249 that error-controlling *procedures* for the random hypothesis case are not as
 250 abundant and studied as the fixed hypothesis case. The mathematical proofs
 251 would typically require some difficult to justify assumptions. Simulation
 252 performances however, do seem promising [10, 9].

253 **3.4 Functional Magnetic Resonance Imaging– Clinical** 254 **Scan**

255 Return again to the neuroscientist from sec 3.3. This time his single patient
 256 is about to enter surgery for the removal of a brain tumor. The patient
 257 will be scanned in the fMRI in order to localize the speech regions, as the
 258 tumor is residing nearby and the surgeon needs to be extra-careful around
 259 these regions. In this clinical case there are different considerations than in
 260 basic scientific research. Type *II* errors are arguably more important than
 261 type *I* errors: underestimating the speech region might cost the patient his
 262 verbal skills; overestimating it, might cost him an extra surgery or a recurring
 263 tumor. None of the error measures presented until now is concerned with false
 264 negatives. Referring to the terminology in Table 1, our neuroscientist would
 265 probably be interested in something like $E(T/(m - R))$, which captures the
 266 sensitivity of the inference. This measure is the False Non Detection Rate
 267 (FNR) [17]. We have not presented this measure yet, as it is concerned with
 268 the *non-detections*. We shall revisit it in the context of power in Section 5.

269 3.5 Genome Wide Association Studies

270 In a typical Genome Wide Association Study (GWAS) the geneticist will
271 record the genetic information of many subjects (genotyping) with the aim of
272 discovering associations between the genotype and the individuals' attributes
273 (phenotype). Assuming a univariate phenotype, the researcher will perform
274 some type of regression between the phenotype and each genetic attribute
275 (titled single nucleotide polymorphism– SNP). With today's technology, the
276 number of SNPs considered in a typical GWAS is hundreds of thousands. To
277 declare an association, a researcher will try to reject the no association null
278 hypothesis between *each* SNP and the phenotype, leading to the simultaneous
279 testing of several hundreds of thousands of hypotheses. Since the researcher
280 does not concern himself with the making of a single mistake, as long as
281 other associations discovered are true, he should choose FDR control or one
282 of its relatives discussed in Section 5. That said, it is also very common in
283 GWAS, to use a p-value threshold of 10^{-7} . This threshold is intended for
284 FWER control, when searching over 500,000 SNPs and using the Bonferonni
285 procedure [8] for FWER control.

286 So FDR or FWER? It is left for the researcher to decide, and it ultimately
287 depends on the implications of declaring false associations.

288 3.6 Cross Validated Example

289 In this example, the researcher is looking for associations between SNPs and
 290 three distinct³ phenotypes. The error measure has already been selected.
 291 The family groupings are unclear. The options being (a) accounting for
 292 errors only within each experiment, leading to three families of hypotheses,
 293 or (b) global error accounting, leading to a single family.

Both approaches have their advantages and disadvantages. By keeping the experiments separate we gain power but the global error rate is no longer α . Yekutieli [32] has approximated, for some cases, that under strategy (a) with α level FDR control within each experiment, then the *global* FDR should actually be close to

$$FDR \approx \alpha \cdot \frac{\text{Total discoveries} + \text{No. of experiments}}{\text{Total discoveries} + 1} \quad (1)$$

294 Eq. 1 captures the intuition that the more experiments performed while
 295 controlling only for errors within the experiment, the global error rate might
 296 inflate.

297 The discussed problem includes three experiments, assumingly looking
 298 for three different phenomena. The fact we discuss *three* experiments, which
 299 were all conducted by the same researcher, is quite arbitrary. Why not control
 300 for the errors performed in the whole of science? Or at least in all genetic
 301 association studies. A proper discussion of this matter requires an unplanned

³ It is actually implicit whether these are distinct phenotypes or not. We have assumed they are distinct, because of the “subject overlap” comment.

302 detour into the philosophy and sociology of science, and is not part of this
 303 guide. We will conclude by remarking that combining errors over different
 304 phenomena is indeed desirable [20], yet rarely performed in practice.

305 **3.7 Imaging Genetics**

306 The field of imaging genetics aims at finding the genetic attributes associates
 307 with phenotypes derived from medical imaging. In a pioneering study, Stein
 308 et al. [25] set out to find the genetic variation associated with local brain
 309 volume, under the paradigm that different genes affect different brain regions.
 310 The data included the genotyping and imaging of $N \approx 700$ individuals. The
 311 genotype of each individual comprises information of $n_G \approx 400K$ SNPs. The
 312 imaging data encodes the relative volume of each subject at $n_B \approx 30K$ voxels.
 313 Testing for association between all $\{SNP\} \times \{voxel\}$ combinations, leads to
 314 $n_G \cdot n_B \approx 12 \times 10^9$ hypotheses. Should they all be considered one family of
 315 hypotheses? Or maybe each SNP (or voxel) is actually a separate family?

316 A researcher might want to infer which gene is associated with which
 317 location. A single family of hypotheses will include all $\{SNP\} \times \{voxel\}$
 318 combinations. FWER control over this family is out of the question. FDR
 319 control means that the researcher is concerned with the proportion of false
 320 associations detected within all of the $\{SNP\} \times \{voxel\}$ associations found.
 321 This seems like a good criterion, except for the fact it requires correcting
 322 for 12×10^9 hypotheses. Our researcher starts considering alternative error
 323 criteria. A natural option might be SNP-wise testing, perhaps using the B-

324 H procedure over all voxels within each SNP. Power is certainly gained as
 325 each family is corrected only for the n_B voxels within it. What about false
 326 findings? Sadly, this approach offers no error control. To see this, consider
 327 the case where there is only one voxel: this amounts to n_G level α hypothesis
 328 tests, which is this initial multiplicity problem.

329 A more justifiable solution might harness the hierarchy of the problem;
 330 that is, by selecting associated SNPs and localizing the association only for
 331 these selected SNPs. Naturally, the multiplicity is alleviated since only se-
 332 lected SNPs will be passed for voxel-wise testing. However, if the same data
 333 is used for selecting the SNPs and then selecting the voxels, an α level FDR
 334 control within selected SNPs will still not guarantee an α level FDR control,
 335 across discovered $\{SNP\} \times \{voxel\}$ associations. This is actually a case of
 336 *selective inference* [2], also referred to by practitioners as “data snooping” or
 337 “double dipping”.

338 Having given it more thought, our researcher decides she wants a method
 339 that has two properties: (a) controlling for the number of falsely discov-
 340 ered SNPs; and (b) controlling for the number of falsely discovered voxels
 341 associated with each discovered SNP.

342 To put it formally, Table 1 needs refinement. Define R and V to be
 343 the number of discovered SNPs and *falsely* discovered SNPS respectively.
 344 Define R_g to be number of voxels declared associated with SNP g , and V_g
 345 accordingly. The desired measure of error has two requirements:

$$E\left(\frac{V}{R}\right) \leq \alpha_1 \text{ and } E\left(\frac{1}{R} \sum_g \frac{V_g}{R_g}\right) \leq \alpha_2 \quad (2)$$

Is there a procedure that controls this type of error? While novel and under active research, there is presently one such procedure⁴. It has two stages. First, the omnibus-stage: testing for an associated SNP by aggregating over voxels within SNP and controlling for the number of SNPs tested. Second, a post-hoc stage: drilling into the selected SNPs searching for associated voxels. The novelty of the procedure is at the second stage, which controls for the number of voxels with a conservative error rate, which accounts for the previous SNP selection stage. The details can be found in [3].

4 Simultaneous versus Selective Inference

Up to this point, we have motivated the choice of error measure by a mere “error accounting”. There is actually another perspective, which can make the choice of the error measure quite obvious once it has been recognized. As put by Cox [11]:

It might be better to talk about the problem of selected comparisons rather than about the problem of multiple comparisons.

Cox’s insight was that making statements on the truthfulness of a subset of selected hypotheses, and making statements about the simultaneous

⁴ With proofs for the independent test-statistics case

363 correctness of a subset of hypotheses is not the same thing. Think in terms
364 of replicability: replicating a combination of phenomena is not the same as
365 replicating each separately. Naturally, the simultaneous truthfulness of the
366 selected hypotheses, entails the truthfulness of each and every one of them.
367 Thus, simultaneous inference is the more ambitious task. In Cox’s words
368 [11]:

369 *The fact that a probability can be calculated for the simultaneous*
370 *correctness of a large number of statements does not usually make*
371 *that probability relevant for the measurement of the uncertainty*
372 *of one of the statements. If we are directly interested in a single*
373 *statement about the vector parameter, the probability of simulta-*
374 *neous correctness would, however, be appropriate. The practical*
375 *usefulness of the multiple comparison techniques then usually lies*
376 *in giving a conservative bound for the effect of selection, rather*
377 *than in giving an “exact” solution.*

378 Armed with the distinction between simultaneous and selective inference,
379 we can relate error measures to inference types.

380 Simultaneity ambitions are only satisfied with FWER control. This is
381 simply because allowing any errors to filtrate, ruins the truthfulness of the
382 combination of claims. FDR is clearly a selective, and not simultaneous state-
383 ment. It can be seen as the average truthfulness of the selected hypotheses.
384 FDX is also selective. It ensures a high proportion of truthful statements
385 within the selected hypotheses.

386 Demonstrating using the examples: The background in Tukey’s psycho-
387 logical exams from Section 3.1 was too vague to determine which inference
388 type is appropriate. Both can be advocated. The same goes for the NAEP
389 example from Section 3.2. Assuming the neuroscientist from the fMRI exam-
390 ple in Section 3.3 cares of the truthfulness of each detected location, and not
391 their particular combination, this is a case of selective inference. A similar
392 consideration holds for the case of associated SNPs in the GWAS example in
393 Section 3.5 and $\{SNP\} \times \{voxel\}$ associations in Section 3.7.

394 **Remark 4.1.** Although not the scope of this manuscript, it is only appro-
395 priate to mention False Coverage Rate adjusted confidence intervals. These
396 intervals estimators, suggested by Benjamini and Yekutieli [6], make the
397 distinction between selective and simultaneous clear. They are not simulta-
398 neous, as they do not offer joint coverage of the selected parameters. They
399 do however offer a high average coverage over the selected parameters.

400 5 Power Considerations

401 The reader might have noticed that the different error measures in Section 2
402 care only of the number of discoveries and false discoveries. Our interest
403 in detection sensitivity is naturally implicit in the procedures researchers
404 employ. Otherwise, never rejecting any null hypothesis, will trivially control
405 all of the error types in Section 2. Power can benefit from (a) knowledge of
406 the proportion of signals in the noise ($1 - m_0/m$) or (b) from an introduction

407 the expected deviations from the null hypotheses.

408 To demonstrate (a), consider two researchers doing the same research.
409 The first, which did some more reading on the topic, knows with complete
410 certainty that there are 10 false null hypotheses (signals) in his 100 hypothe-
411 ses tests. The second, being less thorough, has no access to this knowledge.
412 Naturally, the first can exploit this information. As a trivial example, he will
413 know that more than 10 rejections will certainly contain errors⁵.

414 To demonstrate (b), consider a scenario where the researcher is certain
415 that if an effect exists, it would be of magnitude ± 7 (in some arbitrary
416 scale, say z-values). Performing a one-sample z or t test, while ignoring this
417 belief, will lead the researcher to reject all hypotheses with large (absolute)
418 effects. Particularly, an effect of, say 20, will be considered very extreme, with
419 infinitesimally small p-values. But, when considering the fact that effects
420 are expected to be near 7, the researcher might actually prefer to reject
421 effects near 7 *before* he rejects effects near 20. In statistical terminology,
422 this is simply an underpowered test constructed for the wrong alternative
423 hypothesis.

424 Specifying the expected deviation from the null for each hypothesis tested
425 is no easy task. There exist however, several procedures which use the mul-
426 titude of hypotheses tested in an attempt to empirically characterize the
427 deviations from the null ⁶, and harness this information to gain power. Es-

⁵ This does not mean that the first 10 are necessarily true. Only that more than 10 will certainly contain errors.

⁶ Under mild assumptions regarding the form these deviations might take. Essentially

428 sentially all rely on estimators of the probability of a hypothesis being a true
 429 null given the value of some test statistic z_i . This is the posterior probability
 430 of the null, also named “local fdr” and denoted by $fdr(z_i)$. Details can be
 431 found in [13]. Also see Remark 5.1 on the existence and interpretation of
 432 this probability.

433 This magnitude— the probability of being null given the data— is not
 434 an error rate but rather a test statistic. It is however a rather intuitive
 435 test statistic. So much so that many authors set the rejection criterion to,
 436 say, $fdr(z_i) \leq 0.2$. This lends itself to the interpretation that results with
 437 frequencies smaller than 1 in 5, under the null assumption, are “dangerously
 438 prone to wasting investigators’ resources” [13].

439 Storey [26] establishes a relation between $fdr(z_i)$ and the Marginal FDR
 440 from Section 2.3 so that a researcher opting for the $fdr(z_i) \leq 0.2$ crite-
 441 rion, can receive some sense of how many errors per discovery he will be
 442 doing on average. The relation is data dependent. In the problem analyzed
 443 by Efron [13], rejecting for $fdr(z_i) \leq 0.2$, is approximately equivalent to
 444 marginal FDR control of 0.1. This relation is even more appealing, since
 445 with a growing number of hypotheses being tested, the marginal FDR is a
 446 good approximation of the FDR.

447 Returning to power considerations, and using the notation presented in
 448 Table 1, we can specify many error measures which capture the idea of max-

 assuming that deviations from the null are not uniformly dispersed but rather tend to
 clump together.

449 imal power subject to the false detections being kept at a low level:

$$\min\{FNR \text{ such that } FDR \leq \alpha\} \quad (3)$$

$$\min\{mFNR \text{ such that } mFDR \leq \alpha\} \quad (4)$$

$$\min\{T \text{ such that } V \leq \alpha\} \quad (5)$$

450 The different procedures aimed at satisfying these error functions typi-
 451 cally use the local fdr statistic ($fdr(z_i)$) as a test statistic, but differ in the
 452 heuristics used to compute this statistic [eg. 27, 13, 28], typically, assuming
 453 a large number of hypotheses being tested and independence between test
 454 statistics.

455 The search for procedures with desirable properties under errors mea-
 456 sures such as eq 3 - eq 5, is ongoing. It is an active field of investigation with
 457 beautiful theory, developed either by the Multiple Comparisons Community
 458 or borrowing from statistical decision theory [see 28]. The analysis of the
 459 finite-sampling properties of suggested procedures with respect to some de-
 460 sirable error measures is not an easy task. For a complete treatment, and
 461 state of the art procedures, the reader is referred to [14].

462 **Remark 5.1.** Assigning a *posterior* probability to a hypothesis being null
 463 requires also assigning a *prior* probability to this event. The interpretation
 464 of this probability has raised some controversy [e.g. 22] as it can be seen
 465 as a statement of subjective beliefs, of sampling frequencies or merely as

466 descriptive. Settling this interpretation issue is outside the scope of this
467 manuscript.

468 6 Acknowledgments

469 I wish to thank the many friends and colleagues who contributed of their
470 time and knowledge, to improve clarify this manuscript: Prof. Yoav Ben-
471 jamini, Dr. Daniel Yekutieli, Dr. Jelle Goeman, Dr. Aldo Solari, Kornelius
472 Rohmeyer, David Golan, Neomi Singer and Samuel Cohen.

473 References

- 474 [1] Y. Benjamini. Discovering the false discovery rate. *Journal of the*
475 *Royal Statistical Society. Series B: Statistical Methodology*, 72(4):405–
476 416, 2010.
- 477 [2] Y. Benjamini. Simultaneous and selective inference: Current successes
478 and future challenges. *Biometrical Journal*, 52(6):708721, 2010.
- 479 [3] Y. Benjamini and M. Bogomolov. Adjusting for selection bias in testing
480 multiple families of hypotheses. *Journal of the Royal Statistical Society.*
481 *Series B (accepted)*, 2013.
- 482 [4] Y. Benjamini and H. Braun. John w. tukey’s contributions to multiple
483 comparisons. *The Annals of Statistics*, 30(6):1576–1594, December 2002.

- 484 [5] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a
485 practical and powerful approach to multiple testing. *JOURNAL-ROYAL*
486 *STATISTICAL SOCIETY SERIES B*, 57:289–289, 1995.
- 487 [6] Y. Benjamini and D. Yekutieli. False discovery rate-adjusted multiple
488 confidence intervals for selected parameters. *Journal of the American*
489 *Statistical Association*, 100(469):71–81, 2005.
- 490 [7] F. Bretz, T. Hothorn, and P. Westfall. *Multiple Comparisons Using R*.
491 Chapman and Hall/CRC, 1 edition, July 2010. ISBN 1584885742.
- 492 [8] W.S. Bush and J.H. Moore. Chapter 11: Genome-wide association stud-
493 ies. *PLoS Computational Biology*, 8(12), December 2012.
- 494 [9] J. Chumbley, K. Worsley, G. Flandin, and K. Friston. Topological FDR
495 for neuroimaging. *NeuroImage*, 49(4):3057–3064, 2010.
- 496 [10] J.R. Chumbley and K.J. Friston. False discovery rate revisited: FDR
497 and topological inference using gaussian random fields. *NeuroImage*, 44
498 (1):62–70, January 2009.
- 499 [11] D. R. Cox. A remark on multiple comparison methods. *Technometrics*,
500 7(2):223–224, 1965.
- 501 [12] D. Donoho and J. S. Jin. Higher criticism for detecting sparse het-
502 erogeneous mixtures. *Annals of Statistics*, 32(3):962–994, June 2004.
503 WOS:000221981400005.

- 504 [13] B. Efron. Microarrays, empirical bayes and the two-groups model. *Sta-*
505 *tistical science*, 23(1):1–22, 2008.
- 506 [14] B. Efron. *Large-Scale Inference: Empirical Bayes Methods for Estima-*
507 *tion, Testing, and Prediction*. Cambridge University Press, 1 edition,
508 September 2010. ISBN 0521192498.
- 509 [15] A. Farcomeni. A review of modern multiple hypothesis testing, with
510 particular attention to the false discovery proportion. *Statistical Methods*
511 *in Medical Research*, 17(4):347–388, August 2008.
- 512 [16] Helmut Finner and M. Roters. On the false discovery rate and expected
513 type i errors. *Biometrical Journal*, 43(8):985–1005, 2001.
- 514 [17] C.R. Genovese and L. Wasserman. Operating characteristics and ex-
515 tensions of the false discovery rate procedure. *Journal Of The Royal*
516 *Statistical Society Series B*, 64(3):499–517, 2002.
- 517 [18] C.R. Genovese and L. Wasserman. Exceedance control of the false dis-
518 covery proportion. *Journal of the American Statistical Association*, 101
519 (476):1408–1417, December 2006.
- 520 [19] Jelle J. Goeman and Aldo Solari. Tutorial in biostatistics: multiple
521 hypothesis testing in genomics. *Statistics in Medicine*, 2013.
- 522 [20] John P. A. Ioannidis. Why most published research findings are false.
523 *PLoS Medicine*, 2(8):e124 EP –, August 2005.

- 524 [21] N.A. Lazar. *The Statistical Analysis of Functional MRI Data*. Springer,
525 1 edition, July 2008. ISBN 0387781900.
- 526 [22] Carl N. Morris. Comment: Microarrays, empirical bayes and the two-
527 groups model. *Statistical Science*, 23(1):34–40, February 2008. Article-
528 Type: research-article / Full publication date: Feb., 2008 / Copyright
529 2008 Institute of Mathematical Statistics.
- 530 [23] R Development Core Team. R: A language and environment for statisti-
531 cal computing. <http://www.R-project.org>, 2011. URL [http://www.R-](http://www.R-project.org)
532 [project.org](http://www.R-project.org).
- 533 [24] D. O. Siegmund, N. R. Zhang, and B. Yakir. False discovery rate for
534 scanning statistics. *Biometrika*, 98(4):979–985, December 2011.
- 535 [25] J.L. Stein, X. Hua, S. Lee, A.J. Ho, A.D. Leow, A.W. Toga, A.J.
536 Saykin, L. Shen, T. Foroud, N. Pankratz, M.J Huentelman, D.W. Craig,
537 J.D. Gerber, A.N. Allen, J.J. Corneveaux, B.M. DeChairo, S.G. Potkin,
538 M.W. Weiner, and P.M. Thompson. Voxelwise genome-wide association
539 study (vGWAS). *NeuroImage*, 53(3):1160–1174, November 2010.
- 540 [26] J. D. Storey. The positive false discovery rate: A bayesian interpretation
541 and the q-value. *ANNALS OF STATISTICS*, 31(6):2013–2035, 2003.
- 542 [27] J.D. Storey. A direct approach to false discovery rates. *Journal Of The*
543 *Royal Statistical Society Series B*, 64(3):479–498, 2002.

- 544 [28] W. Sun and T. Cai. Oracle and adaptive compound decision rules for
545 false discovery rate control. *Journal of the American Statistical Association*,
546 102(479):901–912, September 2007.
- 547 [29] M.J. van der Laan, S. Dudoit, and K.S. Pollard. Augmentation pro-
548 cedures for control of the generalized family-wise error rate and tail
549 probabilities for the proportion of false positives. *Statistical applications*
550 *in genetics and molecular biology*, 3:Article15, 2004.
- 551 [30] P.H. Westfall, R.R.D. Tobias, and R.D. Wolfinger. *Multiple Comparisons*
552 *and Multiple Tests Using SAS, Second Edition*. SAS Institute, August
553 2011. ISBN 9781607648857.
- 554 [31] Valerie SL Williams, Lyle V. Jones, and John W. Tukey. Controlling
555 error in multiple comparisons, with examples from state-to-state differ-
556 ences in educational achievement. *Journal of Educational and Behavioral*
557 *Statistics*, 24(1):4269, 1999.
- 558 [32] D. Yekutieli. Hierarchical false discovery RateControlling methodology.
559 *Journal of the American Statistical Association*, 103:309–316, March
560 2008.

561 **A On Your Computer**

562 It does not suffice to choose an error measure in order to perform an analysis.
563 An error-controlling procedure will also have to be chosen, and this is what

564 you should look for in your favorite software. This might be a general purpose
565 statistical suite, or a problem-specific application. In the latter case, we
566 have little to suggest, as domain specific applications typically implement
567 the procedures popularized in that field. In the brain imaging example,
568 popular software include SPM, Brain Voyager, FSL, AFNI. All incorporate
569 the multiplicity control procedure preferred by their authors (thus implicitly,
570 the error measure). In the GWAS example, the same occurs in software such
571 as Plink, PRESTO, PERMORY and others. General purpose statistical
572 software are built for flexibility in the analysis, and thus incorporate more
573 multiplicity control procedures.

574 In the R programming environment [23] the function *p.adjust* in the *stats*
575 package will allow you to perform the most common procedures. The refer-
576 ences in the function documentation are a good starting point for learning
577 about these procedures. For FWER controlling procedures, in particular in
578 the context of linear contrasts in regression models, the *multcomp* package
579 is a good option. For FDR control (and variants) many packages have been
580 written. A good listing of these can be found in Bretz et al. [7] or Korbinian
581 Strimmer's web site: <http://strimmerlab.org/notes/fdr.html>. We also
582 note that, to the best of our knowledge, the hierarchical testing scheme in
583 Section 3.7 has not been implemented. Its implementation would require a
584 new syntax to describe the hypotheses' hierarchy (families) and it is an open
585 challenge.

586 In SAS, multiple testing procedures are incorporated within PROC MIXED

587 and also in PROC MULTTEST. The canonical reference is Westfall et al. [30].

588 In SPSS, multiplicity corrections are typically found as part of the *post-*
589 *hoc* options of the analysis methods.

590 **B Glossary**

591 Some of the terms in the multiple-comparisons literature, have appeared in
592 several other disciplines under different names. To ease the transition, and
593 for completeness, we present a glossary. In this glossary, we use as a reference
594 the statistical nomenclature in Table 1. Also note that we use *rate* for the
595 average of a *ratio* or a *proportion*. The literature is not consistent regarding
596 this convention, so that the terms might be found in use for both purposes.

Table 2: Glossary

Symbol	Names
S	True Positives, Hits, True Discoveries
U	True Negatives
V	False Positives, Type <i>I</i> Errors, False Discoveries, False Alarms
T	False Negatives, Type <i>II</i> Errors, False Non Discoveries, Misses
V/R	False Discovery Proportion (FDP), False Discovery Ratio, False Detection Ratio/Proportion, False Alarm Ratio/Proportion, False Positive Ratio/Proportion, Fall-Out
$E(V/R)$	False Discovery Rate (FDR), False Detection Rate, False Alarm Rate, False Positives Rate
R/m	Accuracy
$E(T/(m - R))$	False Non Discovery Rate (FNR)
$T/(m - R)$	False Non Discovery Ratio
S/R	Positive Predictive Value (PPR), Precision, Hit Ratio
$E(S/R)$	Hit Rate
$U/(m - R)$	Negative Predictive Value
$E(T)/m_1$	Non Discovery Rate
T/m_1	Non Discovery Ratio
S/m_1	True Positive Ratio, Recall, Average Power
U/m_1	Specificity, True Negative Ratio
$E(U)/m_1$	True Negative Rate
$E(S/R)$	True Positive Rate