

Better-Than-Chance Classification for Signal Detection

Jonathan Rosenblatt Roei Gilron Roy Mukamel

August 6, 2016

Abstract

[TODO]

1 Introduction

A common workflow in neuroimaging consists of fitting a classifier, and estimating its predictive accuracy using cross validation. Given that the cross validated accuracy is a random quantity, it is then common to test if the cross validated accuracy is significantly better than chance using a permutation test. Examples in the neuroscientific literature include ???, and especially the recently popularized *multivariate pattern analysis* (MVPA) framework of ?. This practice is also observed in the genetics literature, but to a lesser extent [??].

To fix ideas, we will adhere to a concrete example. In ?, the authors seek to detect brain regions which encode differences between vocal and non-vocal stimuli. Following the MVPA workflow, the localization problem is cast as a supervised learning problem: if the type of the stimulus can be predicted from the spatial activation pattern significantly better than chance, then a region is declared to encode vocal/non-vocal information. We call this an *accuracy test*, a.k.a. *class prediction* in ?, or *pattern discrimination* in ?.

This same signal detection task can be also approached as a two-group multivariate test. Inferring that a region encodes vocal/non-vocal information, is essentially inferring that the spatial distribution of brain activations is different given a vocal/non-vocal stimulus. As put in ?:

... the problem of deciding whether the classifier learned to discriminate the classes can be subsumed into the more general question as to whether there is evidence that the underlying distributions of each class are equal or not.

27 A practitioner may then call upon a two-group location test such as Hotelling's
 28 T^2 [?]. Alternatively, if the size of a brain region is too large compared to
 29 the number of observations, so that the spatial covariance cannot be fully
 30 estimated, then a high dimensional version of Hotelling's test can be called
 31 upon, such as in ? or ?. For brevity, and in contrast to *accuracy tests*, we
 32 will call these two-sample multivariate tests simply *location tests*, also termed
 33 *class comparisons* in ?.

34 At this point, it becomes unclear which is preferable: a location test or an
 35 accuracy test? The former with a heritage dating back to ?, and the latter
 36 being extremely popular, as the 959 citations¹ of ? suggest.

37 The comparison between location and accuracy tests was precisely the
 38 goal of ?, who compared the T^2 location test to the accuracy of *Fisher's*
 39 *linear discriminant analysis* classifier (LDA). By comparing the rates of con-
 40 vergence of the powers to 1, ? concluded that accuracy and location tests are
 41 rate equivalent. Judging by convergence rates alone, not much is (asymptot-
 42 ically) lost by using an accuracy test.

43 Asymptotic relative efficiency measures (ARE) are typically used by statis-
 44 ticians to compare between test statistics with similar rates [?]. The ARE
 45 between Hotelling's T^2 (location) test and Fisher's LDA (accuracy) test in ?
 46 is lower bounded by $\sqrt{2\pi} \approx 2.5$. This means that Fisher's LDA requires at
 47 least 2.5 more samples to achieve the same (asymptotic) power than the T^2
 48 test. In this light, the accuracy test is remarkably inefficient compared to the
 49 location test. For comparison, the t-test is only 1.04 more (asymptotically)
 50 efficient than Wilcoxon's rank-sum test [?], so that an ARE of 2.5 is strong
 51 evidence in favor of the location test.

52 Before discarding accuracy tests as inefficient, we recall that ? ana-
 53 lyzed a *half-sample* holdout. The authors conjectured that a leave-one-out
 54 approach, which makes more efficient use of the data, may have better per-
 55 formance. Also, the analysis in ? is asymptotic. This eschews the discrete
 56 nature of the accuracy statistic, which will be shown to have crucial impact.
 57 Since typical sample sizes in neuroscience are not large, we seek to study
 58 which test is to be preferred in finite samples? Our conclusion will be quite
 59 simple: *location tests almost always have more power than accuracy tests*.

60 The main argument for our statement rests upon the observation that
 61 with typical sample sizes, the accuracy test statistic is highly discrete. Dis-
 62 crete test statistics are known to be conservative [?], since they are insensitive
 63 to mild perturbations of the data, and they cannot exhaust the permissible
 64 false positive rate. The degree of discretization is governed by the number of
 65 samples. In our neuroscience example from ?, the classification is performed

¹GoogleScholar. Accessed on Aug 4, 2016.

66 based on 40 trials, so that the test statistic may assume only 40 possible
 67 values. This number of examples is not unusual if considering this is the
 68 number of subjects, or the number of trial-repeats in an neuroimaging study.

69 The discretization effect is aggravated if the test statistic is highly concen-
 70 trated. For an intuition consider the usage of a the *resubstitution accuracy*
 71 as a test statistic. This statistic simply means that the accuracy is not cross
 72 validated. If the data is high dimensional, the resubstitution accuracy will
 73 be very high due to over fitting [?, Theorem 1]. In an extreme case, the
 74 resubstitution accuracy will be 1 for the observed data, but also for any
 75 permutation. The concentration of resubstitution accuracy near 1, and its
 76 discreteness, render this test completely useless, with a power of 0.

77 To compare the power of accuracy tests and location tests in finite sam-
 78 ples, we perform a simulation study of a battery of test statistics. The main
 79 findings are reported in Section 4, and the intuition for our findings is pro-
 80 vided in Section 6, but first, the problem’s setup.

81 2 Problem setup

82 Let $y \in \mathcal{Y}$ be a class encoding. Let $x \in \mathcal{X}$ be a p dimensional feature vector.
 83 In our vocal/non-vocal example we have $\mathcal{Y} = \{-1, 1\}$ and p , the number of
 84 voxels in a brain region so that $\mathcal{X} = \mathbb{R}^{27}$.

85 Given n pairs of (x_i, y_i) , typically assumed i.i.d., a location test amounts
 86 to testing whether $x|y = 1$ has the the same distribution as $x|y = -1$.
 87 I.e., we test if the multivariate voxel activation pattern has the same dis-
 88 tribution when given a vocal stimulus, as when given a non-vocal stimulus.
 89 An accuracy test amounts to learning a predictive model $\hat{f}(x)$ from some
 90 assumed model class $\hat{f} \in \mathcal{F}$. The prediction accuracy, denoted $T_{\hat{f}}^{acc}$, is de-
 91 fined as the probability of a given classifier \hat{f} of making a correct prediction
 92 $T_{\hat{f}}^{acc} := Prob(\hat{f}(x) = y)$ when given a randomly drawn data point, (x, y) .
 93 A statistically significant “better than chance” estimate of $T_{\hat{f}}^{acc}$ is evidence
 94 that the classes are distinct.

95 2.1 Candidate Tests

96 The design of a permutation test using the prediction accuracy, requires the
 97 following design choices:

- 98 1. How to estimate accuracy?
- 99 2. Is the statistic cross validated or not?

- 100 3. For a K-fold cross validated test statistic: should the data be refolded
101 in each permutation?
- 102 4. Permute labels of features?
- 103 5. For a K-fold cross validated test statistic: should the data folding bal-
104 anced (a.k.a. stratified)?
- 105 6. How many folds?

106 We will now address these questions while bearing in mind that unlike the
107 typical supervised learning setup, we are not interested in an unbiased esti-
108 mate of the prediction error, but rather in the mere detection of a difference
109 between two groups.

110 **How to estimate accuracy?** Given a predictor \hat{f} , a natural test statis-
111 tic is some estimate of its accuracy $T_{\hat{f}}^{acc}$. Complicating matters: very low
112 accuracies, even 0, is evidence that the classes are separated, and we only
113 need to invert the predictions. We can thus consider $|T_{\hat{f}}^{acc} - 0.5|$ as the test
114 statistic. This, however, implies that if the classes are identical, random
115 guessing has 0.5 accuracy. This is not true if the classes are not balanced.
116 The chance level in which case is the prevalence of the dominant class, we
117 denote by \hat{p}_{max} . This suggests the following test statistic $|T_{\hat{f}}^{acc} - \hat{p}_{max}|$. Since
118 we will be aggregating these statistics over random data sets where the dom-
119 inant class may have varying frequencies, it seems appropriate to standard-
120 ize the scale of this statistic. We thus also consider the z-scored accuracy:
121 $|T_{\hat{f}}^{acc} - \hat{p}_{max}| / \sqrt{\hat{p}_{max}(1 - \hat{p}_{max})}$.

122 **Cross validate or not?** Were we interested in an unbiased estimator of
123 the prediction error, there is no question that some independent validation
124 is in order. Since we are merely interested in detecting a difference between
125 classes, a biased error estimate is not an issue provided that bias is consistent
126 over all permutations. The underlying intuition is that if the exact same
127 computation is performed over all permutations, then a permutation test
128 will be “fair”, i.e., will not inflate the false positive rate. We will thus be
129 considering both cross validated accuracies, and resubstitution accuracies as
130 our test statistics, a.k.a. *resubstitution classification*.

131 **Refolding?** The standard practice in neuroimaging is to refold the data
132 after each permutation [?]. This is imperative if permuting labels while
133 aiming at balanced data folds. This is not, however, imperative in general.

134 For simplicity, we will adhere to the standard practice of refolding the data
135 within each permutation.

136 **Permute labels of features?** While seemingly identical, the compound-
137 ing of permutations with data foldings renders these two approaches distinct.
138 As an example, consider balanced (stratified) K-fold cross validation where
139 the initial data folding is balanced. After a label permutation, the original
140 folds will probably not be balanced. If the *features* are permuted, then the
141 labels conserve their original fold assignments, and the original folds are bal-
142 anced after each permutation. Since we only report results while refolding
143 the data in each permutation, then the only difference between permuting
144 labels and permuting features seems to be a computational one. We thus
145 adhere to the more common, albeit computationally less efficient practice of
146 permuting labels.

147 **Balanced folding?** As already implied, a standard practice when cross
148 validating is to constrain the data folds to be balanced (i.e. stratified). This
149 is well justified when aiming at unbiased accuracy estimation. This also
150 simplifies matter when aiming at signal detection, as can be seen from the
151 above discussion of the appropriate test statistic. On the other hand, it
152 may complicate matters, as can be seen from the above discussion on label
153 versus feature permutation. We will report results with both balanced and
154 unbalanced data foldings, only to discover, it does not really matter.

155 **How many folds?** Different authors suggest different rules for the num-
156 ber of folds. We will be varying the number of folds. This will affect the
157 concentration of permutation distribution of the estimated accuracy, which
158 will have a crucial effect on the conservativeness of the accuracy test. Our
159 intuition suggests that since more folds imply a less concentrated estimate,
160 then leave-one-out should be the less conservative, and 2-fold should be the
161 most conservative.

162 The of tests we will be comparing is collected for convenience in Table 1.

Name	Basis	CV	Accuracy	Parameters
Hotelling	Hotelling	—	—	shrink=FALSE
Hotelling.shrink	Hotelling	—	—	shrink=TRUE
lda.CV.1	LDA	TRUE	accuracy	—
lda.CV.2	LDA	TRUE	z-accuracy	—
lda.noCV.1	LDA	FALSE	accuracy	—
lda.noCV.2	LDA	FALSE	z-accuracy	—
sd	SD	—	—	—
svm.CV.1	SVM	TRUE	accuracy	cost=1e1
svm.CV.2	SVM	TRUE	accuracy	cost=1e-1
svm.CV.3	SVM	TRUE	z-accuracy	cost=1e1
svm.CV.4	SVM	TRUE	z-accuracy	cost=1e-1
svm.noCV.1	SVM	FALSE	accuracy	cost=1e1
svm.noCV.2	SVM	FALSE	accuracy	cost=1e-1
svm.noCV.3	SVM	FALSE	z-accuracy	cost=1e1
svm.noCV.4	SVM	FALSE	z-accuracy	cost=1e-1

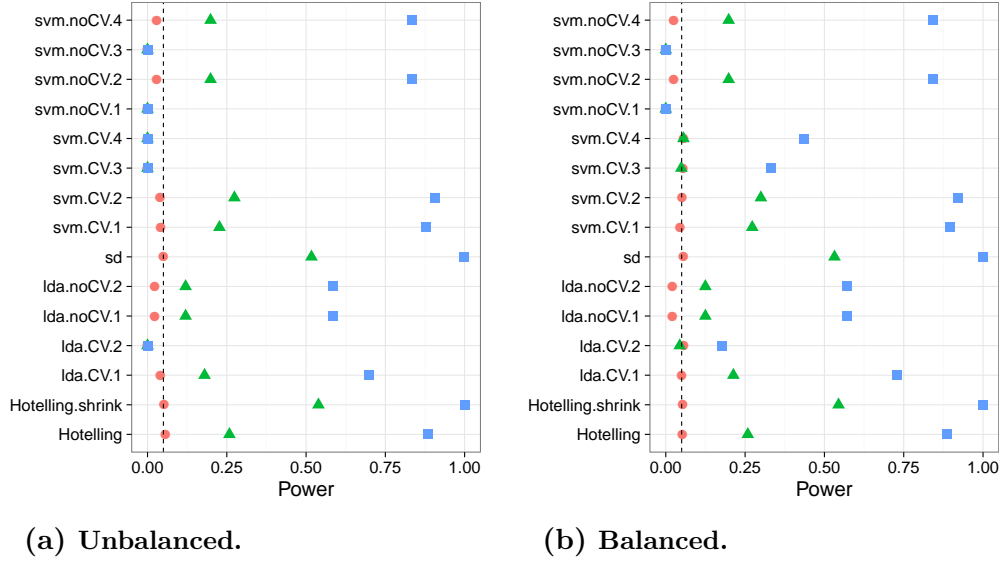
Table 1: This table enumerates the various test statistics we will be studying. Three are location tests: Hotelling, Hotelling.shrink, and sd. *Hotelling* is the classical two-group T^2 statistic. *Hotelling.shrink* is a high dimensional version with the regularized covariance in ?. *sd* is another high dimensional version of the T^2 , from ?. The rest of the tests are variations of the linear SVM, and Fisher’s LDA, with varying accuracy measures, cross validated or not, and varying tuning parameters. For example, *svm.CV.4* is a linear SVM, with *libsvm*’s cost parameter set at 0.1, using the cross validated z-scored accuracy ($|T_{\hat{f}}^{acc} - \hat{p}_{max}|/\sqrt{\hat{p}_{max}(1 - \hat{p}_{max})}$, see Section 2.1). Another example is *lda.noCV.1*, which is Fisher’s LDA, returning the resubstitution accuracy, without cross validation, and without z-scoring.

163

164 3 Controlling the False Positive Rate

165 Figure 1 demonstrates that all of the tests considered conserve the desired
166 0.05 false positive rate, up to varying levels of conservatism. This can be
167 seen from the fact that the probability of rejection is no larger than 0.05 in
168 the absence of any effect, encoded by a red circle. This is true, in particular
169 if: (a) the folds are balanced or not, (b) the tuning parameters of some test
170 statistic are varied, (d) the number of folds is varied. We also observe that
171 the most conservative tests are the resubstitution accuracy measures. We
172 return to this matter in the Discussion.

Figure 1: The power of a permutation test with various test statistics. The power on the x axis. Effect are color and shape coded. The various statistics on the y axis. Their details are given in Table 1. Effects vary over 0 (red circle), 0.25 (green triangle), and 0.5 (blue square). Simulation details in Appendix B. Cross-validation was performed with balanced (stratified) and unbalanced data folding. See sub-captions.



173 4 Power

174 Having established that all of the tests in our battery control the false positive rate, it remains to be seen if they have similar power— especially when
 175 comparing the power of location tests to accuracy tests. From the simulation
 176 results reported in Appendix C we collect the following insights:
 177

- 178 1. Location tests have more power than accuracy tests in all our configurations.
 179
- 180 2. The conservativeness decays as the sample grows (Figures 6a, 6b and
 181 7a), supporting the statement that discretization is responsible for
 182 power loss.
- 183 3. The power may increase or decrease with the number of folds (Figure 3).
 184 [TODO:effect of n.folds.]
- 185 4. The z-scoring of the accuracies was introduced to deal with unbalanced
 186 foldings. If the z-scoring has any effect at all, it merely kills power.
 187 There is really no reason to use it.

- 188 5. Both accuracy and location tests are inappropriate for scale alternatives
189 (Figure 5a). This was to be expected and is reported mostly as a sanity
190 check.
- 191 6. The presence of heavy tails (Figure 5b) may reduce power, but does
192 not quantitatively change results.
- 193 7. Balanced folding typically has no effect. It increased power only for
194 the z-scored statistics (Figure 1). This is surprising given they were
195 precisely designed to deal with the presence of imbalance.
- 196 8. Varying the accuracy test’s tuning parameter, such as the cost (i.e.
197 margins) has no effect on the power of the test.
- 198 9. Correlation between coordinates, mimicking temporal correlation in
199 fMRI data, has no effect on conclusions, since all test statistics account
200 for this correlation (Figure 7b).

201 The major insight from simulations is that the use of accuracy tests for
202 signal detection is underpowered compared to location tests. We now verify
203 this finding on a neuroimaging dataset.

204 5 Neuroimaging Example

205 Figure 2 is an application of both a location and an accuracy test to the data
206 of ?. The authors of ? collected fMRI data while subjects were exposed to the
207 sounds of human speech (vocal), and other non-vocal sounds. Each subject
208 was exposed to 20 sounds of each type, totaling in $n = 40$ trials in each scan.
209 The study was rather large and consisted of about 200 subjects. The data
210 was kindly made available by the authors at the OpenfMRI website².

211 We perform group inference using within-subject permutations using the
212 pipeline of ?, which was also reported in ?. For completeness, the pipeline
213 is described in Appendix A. To demonstrate our point, we compare the *sd*
214 location test with the *svm.cv.1* accuracy test (see Table 1 for the definition
215 of these statistics).

216 In agreement with our simulation results, the location test (*sd*) discovers
217 more brain regions when compared to an accuracy test (*svm.cv.1*). The
218 former discovers 1,232 regions, while the latter only 441, as depicted in
219 Figure 2. We emphasize that both test statistics were compared with the
220 same permutation scheme, and the same error controls, so that any difference
221 in detections is due to their different power.

²<https://openfmri.org/>

222 Having established that accuracy tests are underpowered both in simula-
 223 tion and in application, we wish to identify the conditions under which this
 224 will occur, and discuss implications on the practice of accuracy tests.



Figure 2: Brain regions encoding information discriminating between vocal and non-vocal stimuli. Map reports the centers of 27-voxel sized spherical regions, as discovered by an accuracy test (svm.cv.1), and a location test (sd). svm.cv.1 was computed using 5-fold cross validation, and a cost parameter of 1. Region-wise significance was determined using the permutation scheme of ?, followed by region-wise $FDR \leq 0.05$ control using the Benjamini-Hochberg procedure [?]. Number of permutations equals 400. The location test detect 1,232 regions, and the accuracy test 441, 399 of which are common to both. For the details of the analysis see Appendix A and ?.

225 6 Discussion

226 We have set out to understand which of the tests is more powerful: the
 227 accuracy test or the location test. Using simulations, we have concluded that
 228 the location tests are preferable. We attribute this to several phenomena:
 229 (a) Discretization introduced in finite samples by the accuracy test statistic.
 230 (b) Inefficient use of the data for the validation holdout set. In our high
 231 dimensional setup, we also confirmed that high-dimensional versions of the
 232 T^2 test, such as ? or ? are preferable over the original T^2 .

233 The sensitivity of the power to the number of folds suggests that most of

234 the power is lost due to the discretization and not to the holdout. The degree
235 of discretization is governed by the sample size. For this reason, an asymp-
236 totic analysis such as ? may uncover the holdout inefficiency, but will not
237 uncover the discretization effect. The practical advice for the practitioner, is
238 that for the purpose of signal detection, there is typically a multivariate test
239 (be it a location test or other), that is more powerful than an accuracy test.
240 There is also a good chance that it would be easier to implement, since no
241 validation will be involved.

242 6.1 Neyman-Pearson Classification

243 [TODO]

244 6.2 A good accuracy test

245 In Section 6.5 we discuss cases where an accuracy test cannot replace a
246 location test. For such cases we collect some conclusions from our simulations
247 on the best practices for accuracy tests.

- 248 1. The conservativeness due to discretization decreases with sample size.
- 249 2. Cross-validate. For moderate sample sizes, the power loss due to the
250 holdout inefficiency is smaller than the power loss due to the concen-
251 tration of the resubstitution accuracy.
- 252 3. Permuting features is easier than permuting labels. It allows to preserve
253 balanced folds after a permutation without refolding.
- 254 4. There is no gain in z-scoring the accuracy scores.
- 255 5. Cross validated accuracy with balanced folds has more power than un-
256 balanced folds. We currently have no intuition to offer for this phe-
257 nomenon.
- 258 6. It is unclear what is the effect of the number of folds. More folds in-
259 crease power by reducing the number of holdout samples. On the other
260 hand, it increases the concentration of the accuracy statistic. Com-
261 pounded with the discreteness of the accuracy statistic, this decreases
262 power.
- 263 7. The value of the tuning parameters of a classifier do not matter.

264 6.3 Related Literature

265 ? and ? also looked into a similar problem as we do, namely, what is the
266 preferred accuracy test? They propose a new test they call an *independence*
267 *test*, and demonstrate by simulation that it has more power than other ac-
268 curacy tests, and can deal with non-balanced data sets. We did not include
269 this test in the battery we compared, but we note the following: (a) The
270 independence test of ? relies on a discrete test statistic. This means that in
271 the cases that the accuracy test is called upon for discriminating populations,
272 it will probably be underpowered compared to location tests. (b) In contrast
273 with the underlying motivation of ?’s independence test, we did not find that
274 balancing the data folds is crucial for an accuracy test.

275 6.4 Non-linear predictors

276 It may be argued that accuracy tests permits the separation between classes
277 in high dimensions, such as in *reproducing kernel Hilbert spaces* (RKHS) by
278 using non-linear predictors. This is immaterial since group tests can also be
279 performed in higher dimensions (see ?).

280 6.5 Reservations

281 Some reservations to the generality of our findings are in order. Firstly, not
282 all accuracy tests are concerned with signal detection. Indeed, it is possible
283 that the purpose of the test is not to detect a difference between classes,
284 but to actually test the performance of a particular classifier. Examples
285 include brain decoding for machine interfaces, and clinical diagnosis, where
286 the presence of a medical condition is “predicted” from imaging data. [e.g.
287 ??]

288 Secondly, not all signals are manifested in a shift of the null distribution
289 Our focus on location tests is misleading. Perhaps ?’s *class comparison* is
290 a more appropriate name, in that it does not only imply a shift alternative.
291 Indeed, one may consider signal, i.e. effects, as a change in scale, such as the
292 *spiked covariance* model. In this case, other-than-Hotelling type tests are
293 appropriate [e.g. ?]. Tests have been proposed even when the nature of the
294 difference between populations is left unspecified [e.g. ?]. The fact that in our
295 neuroimaging example (Section 5) some brain regions were detected with the
296 accuracy test, and not the location test, is consistent with this observation.

297 The reservation to the reservation is that the far greater power of the
298 location test, certainly in our example, does serve as an empirical evidence
299 that changes in location are a prevalent phenomenon.

300 6.6 Ease of implementation

301 A very important point is the ease of implementation. The need for cross
302 validation of the accuracy test greatly increases its computational complexity.
303 Moreover, anyone who has actually implemented tests with discrete statistics,
304 will attest they are considerably harder to implement. This is because their
305 unforgiveness to the type of inequality. Indeed, mistakenly replacing a weak
306 inequality with a strong inequality in one's program may considerably change
307 the results. This is not the case for continuous test statistics.

308 6.7 Epilogue

309 Given all the above, we find the popularity of accuracy tests quite puzzling.
310 We believe this is due to a reversal of the inference cascade. Researchers first
311 fit a classifier, and then ask if the classes are any different. Were they to
312 start by asking if classes are any different, and only then try to classify, then
313 location tests would naturally arise as the preferred method. As put by ?:

314 The recent popularity of machine learning has resulted in the ex-
315 tensive teaching and use of prediction in theoretical and applied
316 communities and the relative lack of awareness or popularity of
317 the topic of Neyman-Pearson style hypothesis testing in the com-
318 puter science and related “data science” communities.

319 A Analysis pipeline

320 Here is the analysis pipeline of ? we for the auditory data in ?. Denoting
 321 by $i = 1, \dots, I$ the subject index, $v = 1, \dots, V$ the voxel index, and $s =$
 322 $1, \dots, S$ the permutation index. Since regions³ are centered around a unique
 323 voxel, the voxel index v also serves as a unique region index. Algorithm 1
 324 computes a region-wise test statistic, which is compared to its permutation
 325 null distribution computed by Algorithm 2.

Algorithm 1: Compute a group parametric map.

Data: fMRI scans, and experimental design.
Result: Brain map of group statistics: $\{\bar{T}_v\}_{v=1}^V$

```

1 for  $v \in 1, \dots, V$  do
2   for  $i \in 1, \dots, I$  do
3      $T_{i,v} \leftarrow$  test statistic for subject  $i$  in a region centered at  $v$ .
4    $\bar{T}_v \leftarrow \frac{1}{I} \sum_{i=1}^I T_{i,v}$ .
```

Algorithm 2: Compute a permutation p-value map.

Data: fMRI scans of 20 subjects, experimental design.
Result: Brain map of permutation p-values: $\{p_v\}_{v=1}^V$

```

1 for  $s \in 1, \dots, S$  do
2   permute labels;
3    $\bar{T}_v^s \leftarrow$  parametric map
```

³*searchlight* or *sphere* in the MVPA parlance

328 B Simulation Details

329 The following details are common to all the reported simulations, unless stated
330 otherwise in a figure’s caption. The R code for the simulations can be found
331 in [TODO].

332 Each simulation is based on 4,000 replications. In each replication, we
333 generate n i.i.d. samples from a shift model $\mathbf{x}_i = \mu \mathbf{y}_i^* + \eta_i$. Where $y_i^* = \{0, 1\}$
334 is the class of subject i in dummy coding. Recalling that $y_i = \{-1, 1\}$ is the
335 class in effect coding, then clearly $y_i = 2y_i^* - 1$. The noise is distributed as
336 $\eta_i \sim \mathcal{N}_p(0, \Sigma)$. The sample size $n = 40$. The dimension of the data is $p = 23$.
337 The covariance $\Sigma = I$. Effects, i.e. shifts μ , are equal coordinate p -vectors
338 with coordinates that vary over $\mu \in \{0, 1/4, 1/2\}$.

339 Having generated the data, we compute each of the test statistics in Ta-
340 ble 1. For test statistics that require data folding, we used 8 folds. We then
341 compute a permutation p-value by permuting the class labels, and recomput-
342 ing each test statistic. We perform 400 such permutations. We then reject
343 the $\mu_i = 0$ null hypothesis if the permutation p-value is smaller than 0.05.
344 The reported power is the proportion of replication where the permutation
345 p-value falls below 0.05.

C Simulation Results

Figure 3: Simulation details in Appendix B except the changes in the sub-captions.

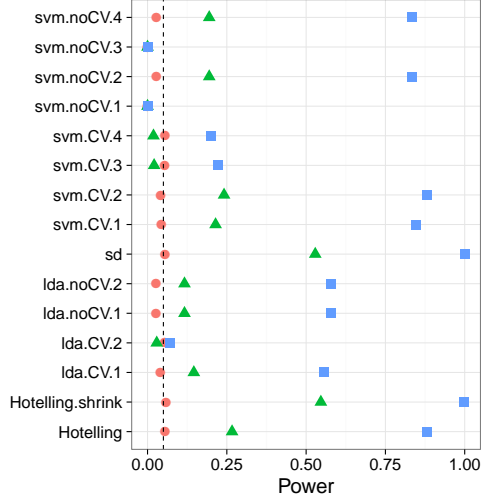


(a) 2-fold cross validation.
Balanced folding.

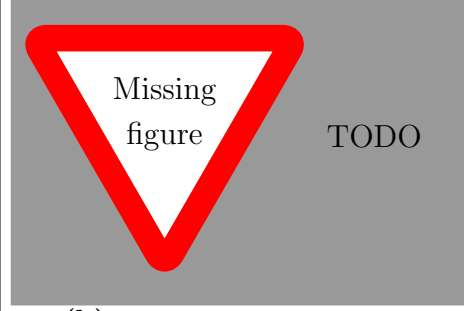


(b) 20-fold cross validation.
Balanced folding

Figure 4: *Simulation details in Appendix B except the changes in the sub-captions.*

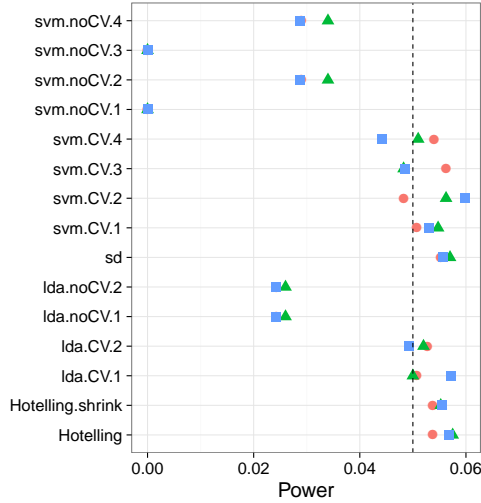


(a) **2-fold** cross validation.
Unbalanced folding.

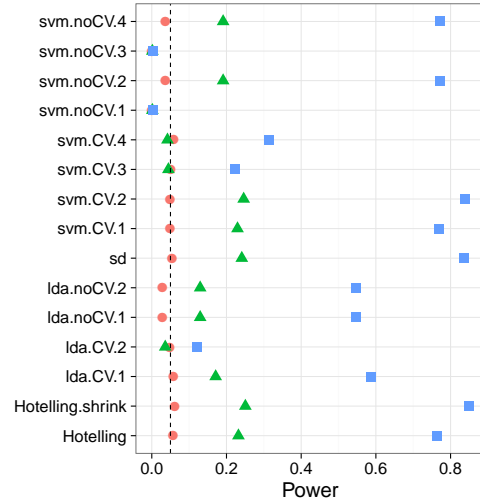


(b) **20-fold** cross validation.
Unbalanced folding.

Figure 5: *Simulation details in Appendix B except the changes in the sub-captions.*

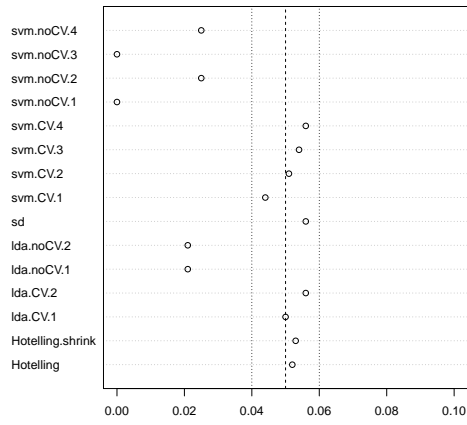


(a) **Scale Change:** $\mathbf{x}_i = \eta_i * \mu^{\mathbf{y}_i^*}$
so that the effect are a scale
change.

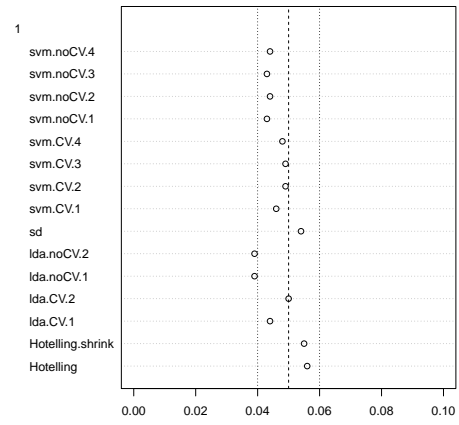


(b) **Heavytailed:** η_i is not
 p -variate Gaussian, but rather
 p -variate t , with $df = 3$.

Figure 6: *Simulation details in Appendix B except the changes in the sub-captions.*

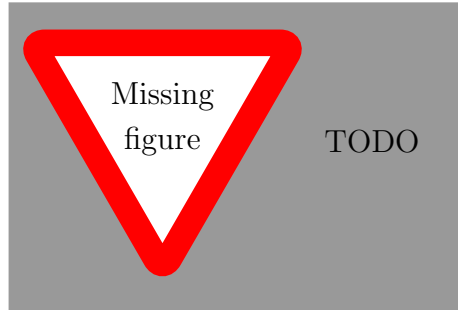


(a) **Low-Dimension:** False positive rates for $n = 40$.



(b) **High-Dimension:** False positive rates for $n = 400$.

Figure 7: Simulation details in Appendix B except the changes in the sub-captions.



(a) **High-Dimension, local alternative:** $n = 400$,
 $\mu \in \frac{\sqrt{40}}{\sqrt{400}} \times \{0, 1/4, 1/2\}$.



(b) **AR(1) dependence:**
 $\Sigma_{k,l} = \rho^{|k-l|}; \rho = 0.8$.