

Surrogate data generation using Deep Knockoffs for nonparametric testing of fMRI time courses

Jonathan Haab
johaab@student.ethz.ch

Supervisor
Dr. M. G. Preti
MIP:lab CIBM EPFL

Mentor
Prof. Dr. D. Iber
CoBi ETHZ

A report presented as part of the
Master in Computational Biology and Bioinformatics
(CBB)

D-BSSE
ETHZ
Switzerland
Spring 2021

Abstract

A precise understanding of the brain requires the study of its structure and activity. Those attributes can be investigated with functional magnetic resonance imaging (fMRI), which has become the major neuroimaging method used for brain mapping thanks to its excellent spatial resolution and its non-invasive nature. In this report, the novel application of *Knockoff Filter* for fMRI data was considered, which would provide an alternative to the phase randomisation technique widely used on such data. The Knockoff methodology provides the considerable advantage of controlling false discovery rate while performing feature selection. We first concentrated our efforts on analysing the fMRI time course surrogates produced using *Deep Knockoffs*, and then employed those surrogates to construct one-sample nonparametric tests at both the individual and group levels. Our results show that this innovative approach while being promising requires more efforts to achieve meaningful outcomes in the specific context of the aforementioned data.

1 Introduction

Through the past decades, functional magnetic resonance imaging (fMRI) has been the soil of substantial discoveries in Neuroscience and stands at this time among the major techniques used in the field. fMRI allows to study the relationship between brain structure and activity and ultimately to build a map of the brain. Such a map is of particular interest to inspect the correlation between structural variations and brain disorders for instance. Interestingly, determining which brain regions are activated upon a particular task execution is similar to a feature (or variable) selection problem, for which the Knockoff framework presents an elegant solution [8].

In this project, blood-oxygen-level dependent (BOLD) contrast fMRI data was analyzed by first fitting a generalized linear model (GLM) (c.f. Fig. 1), making use of the known task paradigm. The GLM was run either once or twice, depending on the type of analysis. The first round allowed to investigate the brain activation at the individual level whereas the second one exposed patterns found at the group level. The GLM produced beta values associated with the degree of activation of a particular region and the significance of those activation scores were then tested.

We used both parametric and nonparametric testing approaches. The former was merely used as a reference point to assess the performance of the later. Indeed, even if parametric tests are effortless to set up and can give high statistical power, they bear some limitations which nonparametric tests can overcome. Since the outcome of nonparametric tests greatly depends on the quality of the surrogates employed, it is of primary interest to develop adroit methods to gen-

erate those surrogates according to an (unknown) null-distribution. Nowadays, the most widely used technique is *Phase randomization* [1], through which time series are transformed into the frequency domain, then modified such that phases are randomized, and transformed back into the time domain hence producing surrogates uncorrelated with the original time courses.

Despite the popularity of phase randomization, we attempted to develop an alternative using the novel Knockoff framework and its latest extension, the Deep Knockoffs [9]. In essence, the Knockoff machine was trained on the available fMRI data in order to produce surrogates following the underlying null-distribution. Those surrogates were then exploited to perform statistical testing. The Knockoff framework developed primarily by R.F.Barber and E.Candès [7] was presented as a clever variable selection procedure allowing control over the false discovery rate (FDR) for more meaningful and replicable results. The Deep Knockoff provides a flexible tool for cases where the covariates distribution is not known at a satisfactory level but a sufficient amount of data is available to learn that distribution.

In this work, we applied the Deep Knockoff methodology to perform nonparametric statistical testing on fMRI time courses and ultimately map the brain functions. We first present some background knowledge about the data at hand, the different testing methods, the Knockoff framework and its newest extension, the Deep Knockoff. Then, we describe the method used to produce brain maps, starting from the data available, exposing the details of the GLM procedure and working our way around the problem of multiple comparisons. We eventually reveal our results and discuss them, before concluding and giving some perspectives for further work.

2 Background

2.1 BOLD fMRI

The BOLD fMRI technique relies on the magnetic properties of deoxyhemoglobin and on the physiological understanding of blood flow [5]. Indeed, arterial and venous blood present different characteristics which provide a way to discriminate them. Also, when a particular brain region is stimulated, more oxygen is required so activated areas can be inferred from oxygen-rich blood flow tracking. Deducing BOLD signal from raw fMRI time courses is a challenging task though, because of high scanner and thermal noise, physiological noise, motion, artifacts and inter-subject variability [6]. Moreover, those time-series data are both temporally and spatially correlated. In this work, the BOLD fMRI time courses were provided by the Center for Biomedical Imaging (CIBM) in Lausanne (Switzerland). The data was collected from 100 subjects across 379 brain regions. Our discussion focuses on one of the seven tasks available, namely the *Motor* task, but the methodology used can easily be applied to the other ones.

2.2 Hypothesis tests

We depend on statistical testing to determine if a certain brain region was activated while the subject performed a specific task. Two families of such testing approaches are commonly defined: parametric and nonparametric. When using parametric tests such as Student's t -test, the data should be continuous and the underlying distribution is assumed to be taking the form of a Gaussian [2]. On the other hand, in the nonparametric setting the data is not assumed to follow any particular distribution. For example, the data distribution could present more skewness or even kurtosis than the Gaussian distribution does. An additional advantage of using nonparametric tests is that they are generally more robust, i.e. less sensitive to the presence of outliers [3]. In this work, our goal was to develop a novel one-sample nonparametric test for fMRI data. The null hypothesis was: "region j is not activated during task condition of interest". The alternative hypothesis was: "region j is either deactivated or activated during that task condition". In a more formal notation:

$$H_0 : \beta_j = 0 \quad H_1 : \beta_j \neq 0 \quad (1)$$

where β_j is the activation of region j for a specific task condition.

2.3 Knockoff framework

The Knockoff method was initially developed for data whose observations follow the classical linear regression model

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{z} \quad (2)$$

where $\mathbf{y} \in \mathbb{R}^n$ is a vector of responses, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a known design matrix, $\beta \in \mathbb{R}^p$ is an unknown vector of coefficient and $\mathbf{z} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ is Gaussian noise [7]. Note that only the case where $n \geq p$ is of interest, since the model would not even be identifiable otherwise. The goal of the Knockoffs is to control the FDR while performing variable selection, i.e. selecting the important features to be included in the model. Note that selecting a feature j corresponds to rejecting the null-hypothesis $H_0 : \beta_j = 0$. The FDR is the proportion of falsely selected features, which can be defined as:

$$\text{FDR} = \mathbb{E} \left[\frac{\#\{j : \beta_j = 0 \text{ and } j \in \hat{S}\}}{\#\{j : j \in \hat{S}\} \vee 1} \right] \quad (3)$$

where $a \vee b = \max\{a, b\}$.

The Knockoff procedure has three steps: **1.** For each feature X_j , a "knockoff" feature \tilde{X}_j is produced in a way that imitates the correlation structure of the original feature. **2.** For each pair of original and knockoff variable, a statistics W_j is calculated which serves as justification for rejecting or not a specific feature. **3.** A threshold based on the data is calculated for the statistics and used to select important features.

A key point for Knockoffs generation is that they must comply with the exchangeability property: [9]

$$(\mathbf{X}, \tilde{\mathbf{X}}) \stackrel{d}{=} (\mathbf{X}, \tilde{\mathbf{X}})_{\text{swap}(j)} \quad \forall j \in \{1, \dots, p\} \quad (4)$$

where $(\mathbf{X}, \tilde{\mathbf{X}})$ is the joint probability distribution of the original features and their Knockoffs, $\stackrel{d}{=}$ indicates equality in distribution and $(\mathbf{X}, \tilde{\mathbf{X}})_{\text{swap}(j)}$ is defined as swapping \mathbf{X}_j with $\tilde{\mathbf{X}}_j$.

The Knockoff procedure has been successfully applied to feature selection problems such as the identification of genetic variation associated with precise phenotype [9] or the selection of cancer biomarkers [11] for example. Even though the original Knockoff method works beautifully well in cases where good models are available to describe the joint distribution of features, an extension relying on deep generative models was developed in order to apply the Knockoff approach to a broader set of application [9].

2.4 Deep Knockoffs

The Deep Knockoffs can be used in cases where we do not have reliable prior knowledge about the distribution of the covariates, allowing the application of Knockoff to a broader range of problems as said earlier. This extension is based on deep generative model and, in our case, moment matching network [9]. Other types of network can be used [10], although those alternatives are not considered in the present work. The details of the Knockoff machine implementation are nicely described by Romano et al. [9] but a noteworthy aspect is that the maximum mean discrepancy (MMD) metric is used to quantify the deviation from exchangeability (as in equ. (4)) and ultimately train the machine to produce valuable Knockoffs. Hence, the machine is trained to minimizing the following expectation:

$$\sum_{j=1}^p \hat{D}_{\text{MMD}} \left[(\mathbf{X}', \tilde{\mathbf{X}}'), (\mathbf{X}'', \tilde{\mathbf{X}}'')_{\text{swap}(j)} \right] \quad (5)$$

where \hat{D}_{MMD} is the empirical estimate of the MMD, $\mathbf{X}', \mathbf{X}'' \in \mathbb{R}^{n/2 \times p}$ are partitions resulting from a random split of the data and $\tilde{\mathbf{X}}', \tilde{\mathbf{X}}''$ the corresponding output of the Knockoff machine. The so trained machine is a model-free Knockoff generator and has the remarkable ability to match higher moments of the underlying null-distribution of the data and, while the previous Knockoff approach required the covariates to follow a known distribution [8].

3 Method

3.1 Generalized linear model (GLM)

GLM is based on the idea that a voxel's time course can be described as a linear combination of a scaled version of the model plus some random noise [5]. As illustrated in Fig. 1a, first-level betas are calculated by fitting the time courses corresponding to the different task conditions on the fMRI time course for a specific region. For the task *Motor*, those task conditions are, in order: movement of left foot, right foot, left hand, right hand and tongue. The instant at which an action was realized are known from the task paradigm. In other

words, we know from the experimental design that the subject was asked to move the right hand from time t_1 to t_2 , then rested before moving the left foot from time t_3 to t_4 . An additional step which was omitted in the schematic representation of the GLM procedure is the convolution of the task condition time courses with the *hemodynamic response function* (HRF). The HRF accounts for the time it took the BOLD signal to rise after a region was activated, giving a physiologically more meaningful guess of what the response should look like. In the end, regions that are significantly correlated with a certain task condition are selected as activated during that condition. Note that the beta value associated with the baseline is of no interest here because it does not bring information linking the task condition and the region activation, therefore only the n first beta values will be designated under the term "first-level betas" in the following.

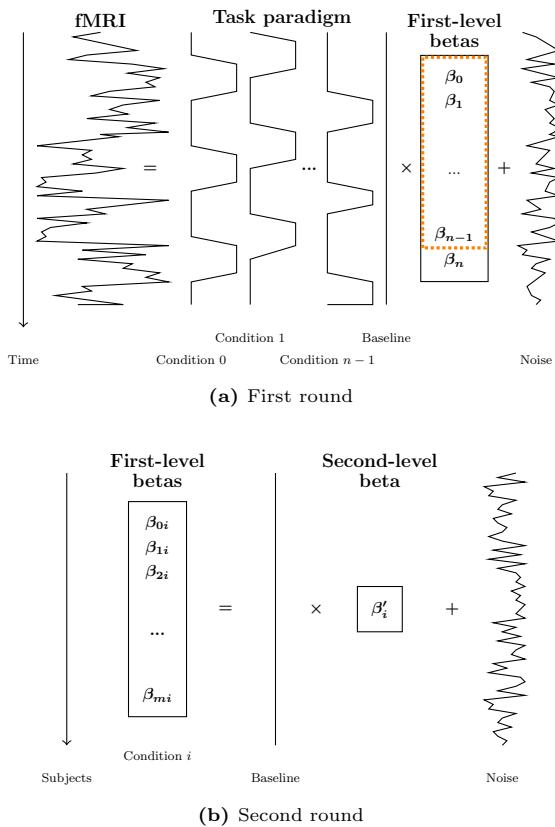


Figure 1: The schemas above present the first and second rounds of GLM procedure. **(a)** In the first round, the different task conditions are fitted on the fMRI time course by minimizing the noise. The beta value associated with the baseline is of no interest for the following steps. The first-level beta values were used for individual analysis. **(b)** In the second round, a baseline is fitted on the first-level beta value of a specific task condition i for every subject. The second-level beta value accounts for the mean activation across the group and was used for group analysis.

Fig. 1b shows the application of GLM for group analysis. In order to identify which regions showed sign of activation across the whole group for a task condition i , The first-level betas from every subject for that task condition are gathered together. A single base-

line is then fitted to those first-level betas, producing a single second-level beta for that region and task condition. If the second-level beta is significantly different from zero, we can conclude that the region in question was (de)activated over the whole group.

3.2 Knockoff surrogates of fMRI time courses

Even though the GLM design is essentially the same as the linear regression model (Equ. (2)) for which the Knockoff framework was first developed, the surrogates were fundamentally different. Since the task conditions used during the experiment are known, no feature selection was needed there. Instead, the Knockoff machine was used to produce surrogate fMRI time courses $\tilde{\mathbf{y}}$ rather than surrogate task condition time courses $\tilde{\mathbf{X}}$. Hence, the question really asked during our nonparametric testing is whether or not region i is significantly more activated in the empirical observation than in its Knockoff copies. The production of Knockoffs showing no substantial sign of activation whatsoever was, and still is, a real challenge.

3.3 Uncorrected and Corrected non-parametric test

In the uncorrected testing approach, a null-distribution was built for each brain region and the test of every region was done independently. Following the method described by Nichols et al. [4], the smallest p -value attainable is $1/N$, where N is the number of surrogates, such that for $\alpha = 0.05$ at least 19 surrogates are needed. In this work, we used $N = 100$ surrogates to have a finer thresholding process. The uncorrected test for individual analysis was performed as follow: **1.** For each surrogate $i = 1, \dots, N$ time course, first-level betas were computed via GLM. **2.** The lower and upper thresholds were taken as the $c + 1$ smallest and largest first-level betas, respectively, where $c = \lfloor \alpha N/2 \rfloor$. **3.** Each first-level beta value resulting from the empirical fMRI time course of a specific subject was tested against the surrogate beta distribution, independently for each region and task condition. If the empirical first-level beta value for that region is smaller (larger) than the $c + 1$ smallest (largest) surrogate beta for that same region, then the region is labeled as deactivated (activated) for the corresponding task condition. Hence, a two-sided test with significant level α was executed.

In the corrected testing approach, a null-distribution was built from minimum (T^{\min}) and maximum (T^{\max}) statistics over the whole surrogate brain. Instead of computing different thresholds for every region, only the minimum and maximum first-level beta surrogates across every region were picked, resulting in N minimum and N maximum statistics. The minimum statistics distribution was only used to identify deactivated regions whereas maximum statistics distribution was for activated regions. The corrected test for individual analysis was performed as follow: **1.** For each surrogate $i = 1, \dots, N$ time course, first-level betas

were computed via GLM. 2. For each surrogate, the minimum and maximum first-level beta value found across the 379 regions were picked, denoted T_i^{\min} and T_i^{\max} respectively 3. The lower threshold was defined as the $c + 1$ smallest beta found in the collection of minimum statistics T^{\min} , where $c = \lfloor \alpha N/2 \rfloor$ as before. Symmetrically, the upper threshold was defined as the $c + 1$ largest beta found in the collection of maximum statistics T^{\max} . Note that the test is still done independently for every task condition, but the minimum and maximum statistics distributions do not vary with the region tested anymore. Also, the use of minimum and maximum statistics make the test more conservative. In brief, two one-sided tests were executed, each with a significant level of $\alpha/2$. For the group analysis, the approaches were similar to the one described above with the difference that second-level beta values were used instead of the first-level betas.

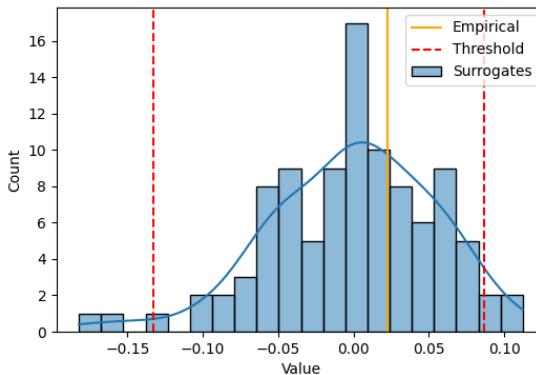


Figure 2: The histogram above shows the distribution of surrogate first-level betas for subject 1, region 100 and task condition 0. The dashed lines highlight the threshold values used for uncorrected nonparametric testing. The orange line highlight the empirical first-level beta, for which the null hypothesis was not rejected in this case (no sign of significant activation of region 100 during task condition 0).

4 Results

4.1 Individual analysis

Individual analysis of fMRI time courses was a natural way to start investigating the performance of Deep Knockoffs surrogates. We looked at the surrogate first-level betas distribution which were used for uncorrected nonparametric testing. Sometimes, the distribution was centered around zero, as was the case for region 100 of Subject 1 during task condition 0 (c.f. Fig. 2), but we found out that for certain pick of regions and task conditions, the whole distribution was shifted toward positive or negative values, as shown in Fig. 3. In those cases, both thresholds are on the same side of the origin, allowing for negligible values to be rejected, thus causing the associated regions to wrongly appear as significantly (de)activates.

Next, the correlation of the surrogate time courses

was inspected. We found that every region of the surrogate brains reached serious correlation with the empirical fMRI time course the Knockoff machine was trained on, and most of them after only two epochs (c.f. Fig. 4).

Finally, we examined the minimum and maximum statistics used for corrected nonparametric testing. As show in Fig. 5, no negligible first-level betas were selected.

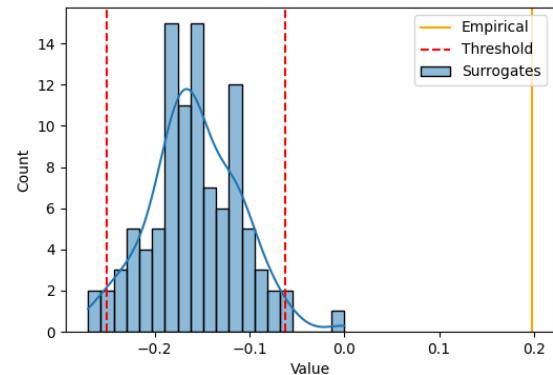


Figure 3: The histogram above shows the distribution of surrogate first-level betas for subject 1, region 100 and task condition 2. The dashed lines highlight the threshold values used for uncorrected nonparametric testing. The orange line highlight the empirical first-level beta, for which the null hypothesis was rejected in this case (significant activation of region 100 during task condition 2). Note that even an empirical first-level beta of value 0 would have been rejected, which obviously represents an erroneous outcome.

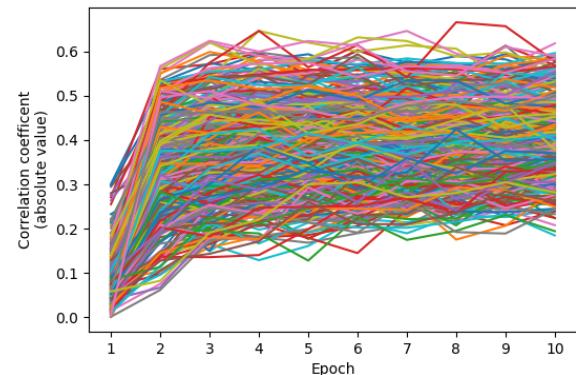


Figure 4: The plot above shows the increase in correlation between the surrogates and the empirical fMRI time course the Knockoff machine was trained on against the number of training epochs. Each one of the 379 lines corresponds to the average over every surrogate correlation coefficient at a specific region.

4.2 Group analysis

The surrogate second-level betas distribution followed the same behavior as shown in Fig. 2 - 5 for the first-level ones. We compared the regions selected by the Knockoff approach with the classical parametric way

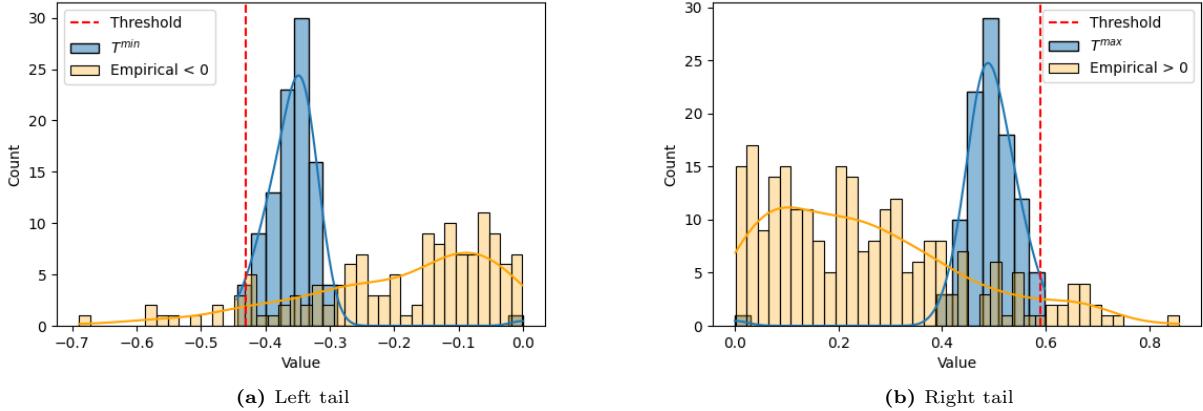


Figure 5: The histograms shown above present the distribution of minimum and maximum statistics, T^{min} and T^{max} , for Subject 1 during task condition 0, in blue. The distribution of the empirical first-level betas is shown in orange. All the empirical betas on the left side of the left tail threshold **(a)** were labeled as significantly deactivated whereas all the empirical betas on the right side of the right tail threshold **(b)** were labeled as activated.

(c.f. Fig. 6) and found that the preeminent region was the same for all methods. The nonparametric techniques seemed to catch more regions than their parametric equivalents though.

5 Discussion

The fact that the Knockoff surrogates got highly correlated with the empirical fMRI time course, as illustrated in Fig. 4, represents a serious problem. The goal of the Knockoff machine was to break the link between the task paradigm and the surrogate time courses. This ambition failed, and was shown to do so by observing the (absolute) correlation coefficient increasing as the number of epochs did. Indeed, if the surrogates were correlated with the empirical time courses, they were also correlated with the initial task paradigm since the empirical time courses were expected to be correlated with the paradigm. The expected output was the perfect opposite: the machine should produce surrogates as different as possible from the example time course used for training, so that the correlation coefficient should decrease as we train the machine. This behavior was unforeseen and resulted in shifted beta surrogates distribution as presented in Fig. 3, causing the uncorrected nonparametric testing to give erroneous output. This effect is most likely the consequence of poor training. Even though the initial data seemed big enough to achieve a satisfying machine training, when looked at more carefully it was noticed that each region of the surrogate brain had only one example at hand to train on. Because each region of the surrogate brain learned independently, it ended up having only one time course from the corresponding region of the subject picked. To solve this issue, we recommend augmenting the training data by using some kind of permutation technique, as circular periodic shift for instance, on the initial empirical time courses. The idea would be to preserve the spatial relationship while shuffling the different time

points. This step would effectively break the temporal dependency of the training data on the task paradigm, so that the surrogates produced would come out uncorrelated with the paradigm.

Aside of the problem related to the machine training, the nonparametric testing using Deep Knockoffs seemed to provide more nuanced output than the parametric alternatives, especially for the corrected tests where the Bonferroni correction might be too conservative (c.f. Fig. 6). The fact that the most activated region was the same for every approach gives good confidence in the novel application of Deep Knockoffs for fMRI data analysis, despite the need for some more fine tuning of the training process. An additional valuable sanity check would be to test white noise and control that the proportion of false positive stays below α , i.e. 5% in our case.

6 Conclusion

In conclusion, we applied the Knockoff filter to implement individual and group analysis via statistical activation mapping, which outputs were comparable to, if not better than, those of parametric testing. We showed that the ability of Deep Knockoffs to preserve the original data characteristics comes at the detriment of high correlation between the surrogates and the training data, a fact which has remained unnoticed until now. Albeit this novel approach appears to be promising, it would require more in-depth investigations of the data the Neural Network is trained on to produce the Deep Knockoffs. In particular, further work should be conducted to explore the usage of circular periodic shift for pre-processing of the empirical fMRI time courses used for training. This small extension might be an neat step towards more meaningful surrogates.

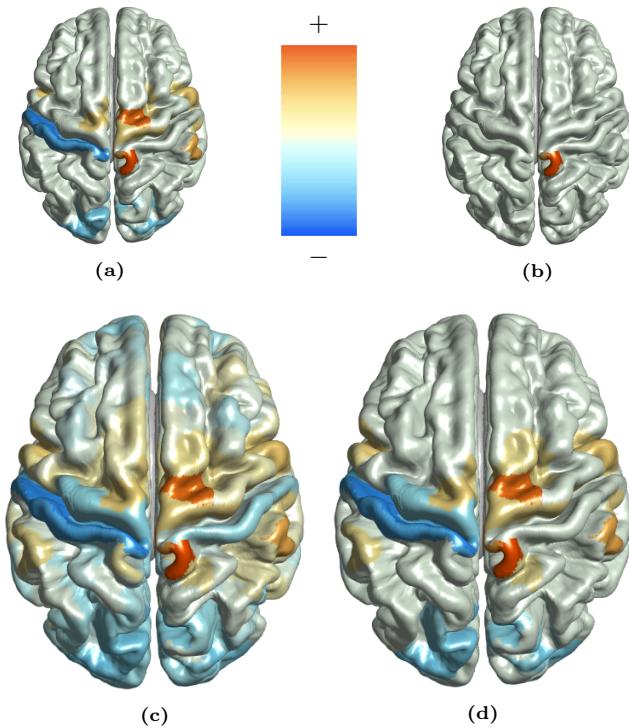


Figure 6: The brain plots above show the (de)activated regions during task condition 0 across the whole group of 100 subjects. Each brain plot results from a different analysis: (a) uncorrected parametric testing, (b) parametric testing using Bonferroni correction, (c) uncorrected nonparametric testing and (d) nonparametric testing using Min/Max correction. We see that the highest activated region is preserved across all approaches, but the nonparametric ones seem to give more nuances than their parametric equivalent.

References

1. Theiler, J., Eubank, S., Longtin, A., Galdrikian, B. & Doyne Farmer, J. Testing for nonlinearity in time series: the method of surrogate data. *Physica D: Nonlinear Phenomena* **58**, 77–94. ISSN: 01672789 (Sept. 1992).
2. Holmes, A. P., Blair, R. C., Watson, J. D. & Ford, I. Nonparametric analysis of statistic images from functional mapping experiments. *Journal of Cerebral Blood Flow and Metabolism* **16**, 7–22. ISSN: 0271678X. <https://pubmed.ncbi.nlm.nih.gov/8530558/> (1996).
3. Douglas Ward B. *Nonparametric Statistical Analysis of FMRI Data* tech. rep. (Biophysics Research Institute, Medical College of Wisconsin, July 1997). <https://afni.nimh.nih.gov/pub/dist/doc/manual/Nonparametric.pdf>.
4. Nichols, T. E. & Holmes, A. P. Nonparametric permutation tests for functional neuroimaging: A primer with examples. *Human Brain Mapping* **15**, 1–25. ISSN: 1065-9471. <http://doi.wiley.com/10.1002/hbm.1058> (Jan. 2002).
5. Buxton, R. B. The physics of functional magnetic resonance imaging (fMRI). *Reports on Progress in Physics* **76**, 096601. ISSN: 00344885. <https://iopscience.iop.org/article/10.1088/0034-4885/76/9/096601%20https://iopscience.iop.org/article/10.1088/0034-4885/76/9/096601/meta> (Sept. 2013).
6. Greve, D. N., Brown, G. G., Mueller, B. A., Glover, G. & Liu, T. T. A Survey of the Sources of Noise in fMRI. *Psychometrika* **78**, 396–416. ISSN: 00333123. <https://link.springer.com/article/10.1007/s11336-012-9294-0> (July 2013).
7. Barber, R. F. & Candès, E. J. in *The Annals of Statistics* **5**, 2055–2085 (Institute of Mathematical Statistics, 2015). https://www.jstor.org/stable/43818570?seq=1#metadata_info_contents.
8. Candes, E., Fan, Y., Janson, L. & Lv, J. Panning for Gold: Model-X Knockoffs for High-dimensional Controlled Variable Selection. *Journal of the Royal Statistical Society. Series B: Statistical Methodology* **80**, 551–577. <http://arxiv.org/abs/1610.02351> (Oct. 2016).
9. Romano, Y., Sesia, M. & Candès, E. J. Deep Knockoffs. *Journal of the American Statistical Association* **115**, 1861–1872. <http://arxiv.org/abs/1811.06687%20http://dx.doi.org/10.1080/01621459.2019.1660174> (Nov. 2018).
10. Jordon, J., Yoon, J. & Van Der Schaar, M. *Knockoff GAN: Generating Knockoffs for Feature Selection using Generative Adversarial Networks* tech. rep. (2019). <https://openreview.net/pdf?id=ByeZ5jC5YQ>.
11. Shen, A., Fu, H., He, K. & Jiang, H. False Discovery Rate Control in Cancer Biomarker Selection Using Knockoffs. *Cancers* **11**, 744. ISSN: 2072-6694. <https://www.mdpi.com/2072-6694/11/6/744> (May 2019).

Supplement

A Code

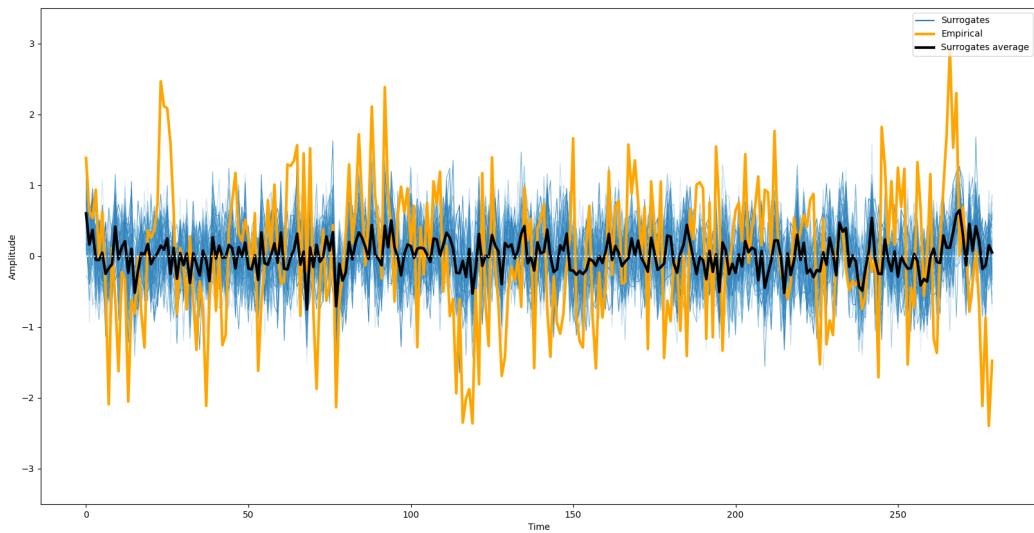
The code used for this project is available on the following GitHub repository:

- https://github.com/johaab/Deep_Knockoffs

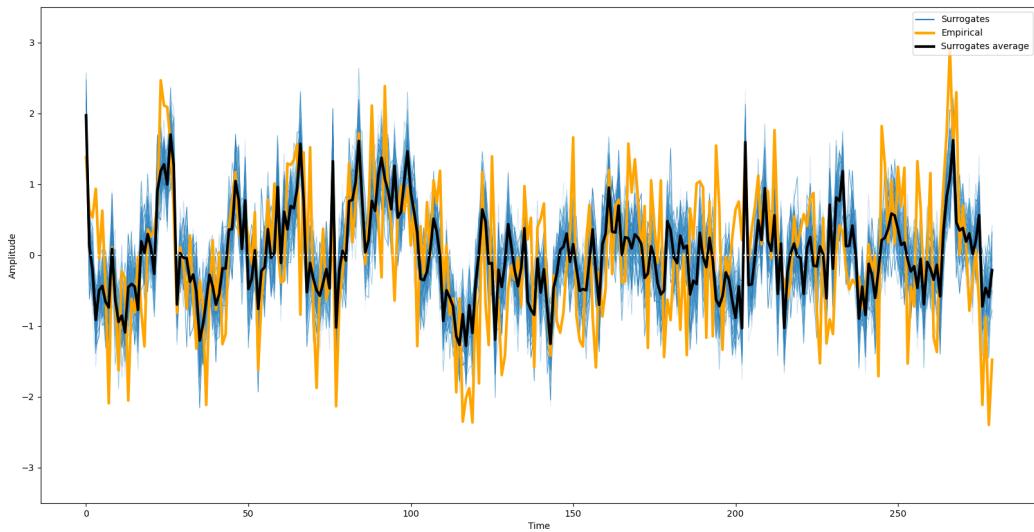
which was merely an extension of the work from Alec Flowers, Alexander Glavackij and Janet van der Graaf:

- <https://gitlab.com/aglavac/machine-learning-cs433-p2>
- itself based on the original Deep Knockoffs implementation:
- <https://github.com/msesia/deepknockoffs>

B Individual analysis



(a) 1 epochs



(b) 10 epochs

Figure B.1: The plots above illustrate the learning process of the Knockoff machine. The empirical fMRI time course which was fed into the Neural Network is shown in orange, while the surrogate time courses produced are shown in the background (blue). The average over the surrogates is shown in black. (a) After the first epoch, the surrogates did not show any pattern. (b) We observe that after only 10 epochs, the surrogates seemed to fit the empirical data very closely. Unfortunately, this behavior presents a problem since the surrogates were expected to be decorrelated with the initial task paradigm.

C Group analysis

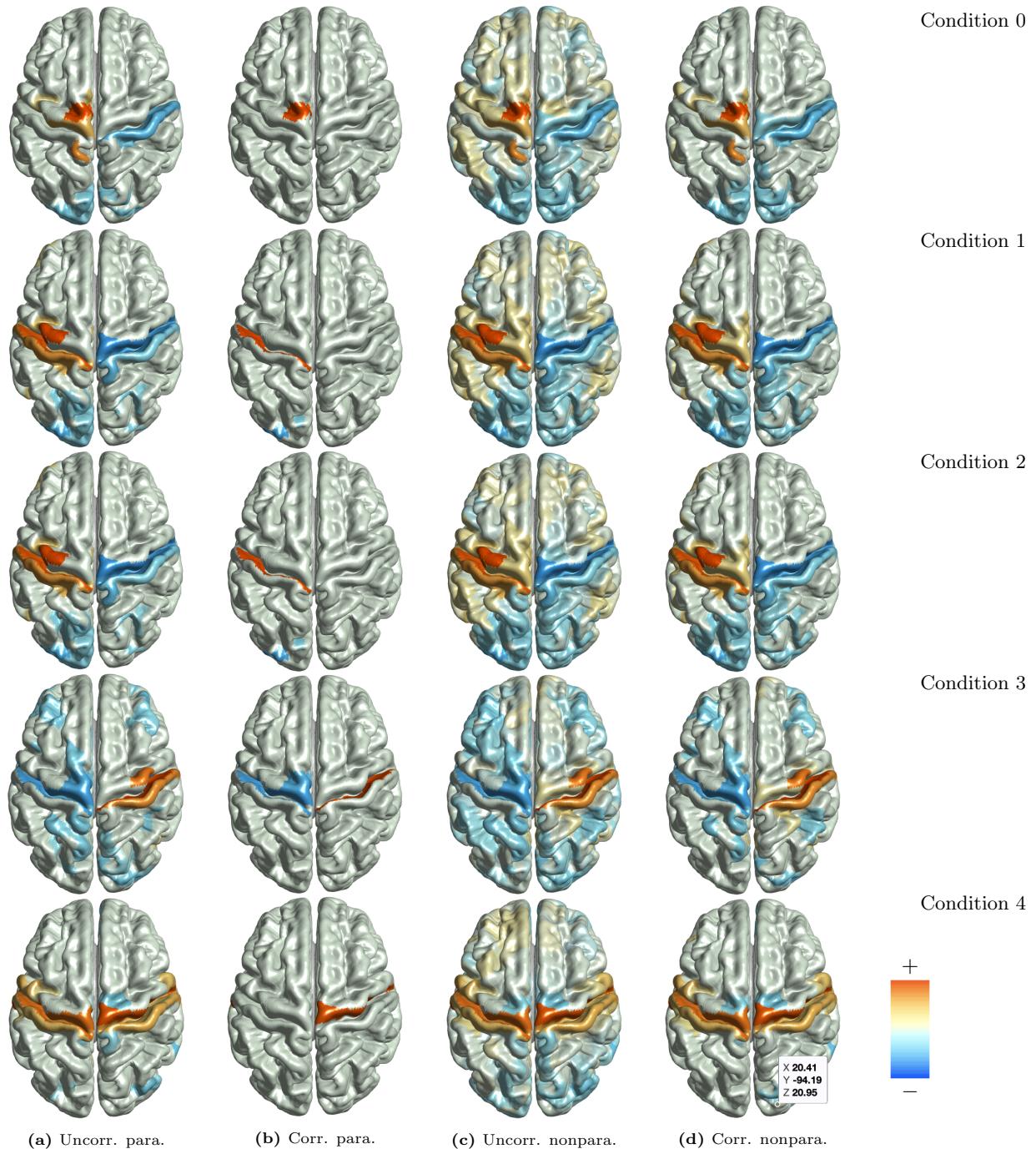


Figure C.1: The brain plots above show the (de)activated regions during task condition 1-4 across the whole group of 100 subjects. Each brain plot results from a different analysis: (a) uncorrected parametric testing, (b) parametric testing using Bonferroni correction, (c) uncorrected nonparametric testing and (d) nonparametric testing using Min/Max correction. We see that the highest activated region is preserved across all approaches, but the nonparametric ones seem to give more nuances than their parametric equivalent.