

Optimising analysis choices for multivariate decoding: creating pseudotrials using trial averaging and resampling

C. L. Scrivener^{1,2}, T. Grootswagers³, A. Woolgar²

¹ Department of Psychology, School of Philosophy, Psychology and Language Sciences, University of Edinburgh, Edinburgh, UK

²MRC Cognition and Brain Sciences Unit, University of Cambridge, Cambridge, UK

³ The MARCS Institute for Brain, Behaviour and Development, Western Sydney University, NSW, Australia

Abstract

Multivariate pattern analysis (MVPA) is a popular technique that can distinguish between condition-specific patterns of activation. Applied to neuroimaging data, MVPA decoding for inference uses above chance decoding to identify statistically reliable condition-specific information in neuroimaging data which may be missed by univariate methods. However, several analysis choices influence decoding success, and the combined effects of these choices have not been fully evaluated. We systematically assessed the influence of trial averaging and resampling on decoding accuracy and subsequent statistical outcome on simulated data. Although the optimal parameters varied with the classifier and cross-validation approach used, we found that modest trial averaging using roughly 5-10% of the total number of trials per condition improved accuracy and associated t-statistics. In addition, a resampling value of 2 could improve t-statistics and classification performance, but was not always necessary. We provide code to allow researchers to optimise analyses for the parameters of their data.

Key words: Multivariate pattern analysis, decoding, pseudotrials.

Introduction

The last decade has seen an explosion in the popularity of multivariate pattern analysis (MVPA) for neuroimaging data. By identifying condition-specific patterns of activation, MVPA can reveal the evolution of information processing over time and/or space and is sensitive to information missed by univariate methods (e.g., Grootswagers, 2017; Pereira et al., 2009; Haynes & Rees, 2006). Typically, cognitive neuroscience experiments employ MVPA techniques, such as linear classification (‘decoding’), to make inferences about the type of information decodable from neuroimaging data, and to characterise these ‘neural representations’ in terms of when and where they can be decoded, whether they generalise between conditions, and if they change with experimental manipulations. Decoding for inference typically compares decoding metrics (e.g., classification accuracy) between conditions or to chance, drawing inference about neural processing from statistically reliable condition-specific information in neuroimaging data. This differs from decoding for prediction which aims to maximise classification performance (Hebart & Baker, 2018).

There are many ways to run a decoding for inference analysis, and multiple decisions are likely to influence decoding success. One analysis option is the creation of ‘pseudotrials’, or ‘supertrials’, by averaging data from subsets of trials together before performing MVPA. Previous results demonstrate that pseudotrial averaging can lead to an increase in classification accuracy, compared to single-trial decoding (e.g., Adam et al., 2020; Tuckute et al., 2019; Hebart et al., 2018; Grootswagers et al., 2017; Isik et al., 2014). However, too much averaging can be detrimental, as increasing the number of trials per pseudotrial can also increase the between-subject variance, which in turn affects the statistical outcome. A second decision is the cross-validation procedure used to evaluate the generalisation of classification across subsets of the data (Bishop & Nasrabadi, 2006). In a leave-one-trial-out procedure, the number of cross-validation steps is equal to the number of trials available per exemplar. At each step, the classifier is trained on all but one of the trials, which is then used to test the classification. A leave-one-pseudotrial-out method uses the same logic, but is based on averaged subsets of the original trials. Another option is to divide the data into larger chunks or blocks across which to train and test the classifier. For example, in a study with 90 trials per exemplar, 10-fold cross-validation would split the data into 9 sets of 10 trials to use iteratively for training and testing. This means that at each iteration, the classifier is trained on fewer data points than for leave-one-trial-out, but previous work has shown that this may provide more stable estimates with lower variance across decoding accuracies (Varoquaux et al., 2017). When combined with trial averaging, pseudotrials can be created either by averaging all of the trials within a chunk or by grouping them into smaller subsets of trials. Therefore, at least two interacting parameters seem likely to affect results: the number of trials averaged together before classification, and the number of folds into which the data are split for cross-validation.

Thus far the choice of these parameters has been largely arbitrary, resulting in a wide variety of trial averaging and cross-validation approaches in the literature. To name a few examples, in Bae and Luck (2018) and Foster et al. (2017), the available trials per condition were randomly divided into 3 chunks before being averaged and decoded using a 3-fold cross-validation. In Isik et al. (2014), groups of 10 trials were averaged before using a 5-fold cross-validation, and in Duncan et al. (2023), groups of 3 trials were averaged before using a 10-fold cross-validation. Goddard et al. (2022) averaged over 16 trials, before using an 8-fold cross-validation, and Petit et al. (2023) used the median of 5 trials and a leave-one-pseudotrial-out

cross-validation with 22 folds. Given the variation across experiments, it can be difficult to make decoding analysis decisions regarding a new set of data.

In an evaluation of different decoding approaches, Grootswagers et al. (2017) compared decoding accuracies when averaging together 4, 8, 16, and 32 trials (equivalent to creating 8, 4, 2, and 1 pseudotrials) from an experiment of 32 trials per condition. A leave-one-pseudotrial-out cross-validation scheme was used, meaning that the number of cross-validation steps was determined by the number of pseudotrials. All averaging procedures increased decoding accuracies compared to no trial averaging, but the least amount of averaging (4 trials, 12.5%) provided the best trade-off between signal-to-noise and number of pseudotrials. In addition, they found similar decoding accuracies for 10-fold and leave-one-trial-out cross-validation methods. They reported lower decoding accuracies using a 2-fold cross-validation procedure, though others have argued that this approach may be associated with higher statistical power overall (Valente et al., 2021). In a different implementation, Adam et al. (2020) compared trial averaging results using a 3-fold classification that was independent from the number of pseudotrials created. They averaged groups of trials within each fold, ranging from 5 to 25, and found a significant increase in the average decoding accuracy with more averaging. However, increasing the number of trials per pseudotrial also increased the between-subject variance, and they concluded that the average of 10 trials was optimal given the number of trials available. Thus, while there is clearly a trade-off between providing a classifier with fewer less noisy trials (more averaging) or more noisy trials (less averaging), it is not yet clear how to optimise this decision. In particular, it unclear whether and how the optimal amount of averaging depends on other factors such as number of trials, choice of classifier, effect size, etc.

The aim of the current work was to inform future decoding studies by systematically assessing the influence of a range of parameters on decoding accuracy and subsequent statistical outcome, using data simulated in CosMoMvPA (Oosterhof, 2016). We varied several parameters across simulations, including the number of trials per pseudotrial, cross-validation procedure, number of trials per condition and the size of the underlying effect, and assessed the data for three different linear classifiers. In all cases, we repeated the averaging procedure multiple times per ‘subject’ to prevent the results relying on one specific division of the trials (Goddard et al., 2022). Decoding results varied with the classifier used, the cross-validation approach, and the number of the trials averaged to create each pseudotrial. In addition, we evaluated the influence of random sampling with replacement (‘resampling’) on decoding accuracy, which has not yet been comprehensively assessed. We hypothesised that this would be beneficial when creating pseudotrials as it increases the number of samples available for classification, which would normally be diminished by trial averaging. Although the optimal parameters varied with classifier and cross-validation approach, we found that using roughly 5-10% of the total number of trials per condition was optimal for creating pseudotrials. In addition, a resampling value of 2 could improve t-statistics and reduce the impact of the number of trials per pseudotrial on classification performance.

Methods

We used CoSMoMvPA (Oosterhof et al., 2016) to simulate multiple datasets, each with two experimental conditions, and examined the influence of a) the number of trials used per pseudotrial, b) the number of times each trial was resampled, c) the simulated class distance,

and d) the choice of classifier. We chose three popular classifiers that are implemented in CoSMoMVPA: linear support vector machine (libSVM), linear discriminant analysis (LDA), and Gaussian Naïve Bayes. For experiment 1, we simulated data with 700 features (reflecting channels/voxels) from 100 ‘subjects’ with 2 conditions and 90 trials per condition. To check whether the results were dependent on the specific number of trials per condition, we also simulated a second experiment with only 45 trials per condition. We also manipulated the multivariate class distance using built-in CoSMoMVPA functions: individual values were drawn from a normal distribution ($sd=1$) with the amount of separation between class-means for each feature defined as $class_mean = class_distance/\sqrt{\log(700*ntrials)}$ with $ntrials$ either 90 or 45 (for the 2 experiments), and used 3 values for $class_distance$ (0, 0.1, and 0.2). Results for a smaller number of participants (50) and fewer features (6) can be found in the supplementary material. Simulation scripts are available on the Open Science Framework, <https://osf.io/hjf75/>.

We ran both a ‘chunking’ procedure and a ‘leave-one-pseudotrial-out’ procedure as both approaches are commonly used. For the ‘chunking’ procedure, pseudotrials were created separately within 3 ‘blocks’ of trials, and a 3-fold cross-validation method was always used (additional results for a 10-fold cross-validation procedure can be found in the supplementary material). For the ‘leave-one-pseudotrial-out’ method, the number of cross-validation steps was determined by the number of trials per pseudotrial (the more trial averaging and resampling, the fewer possible steps). In experiment 1 (90 trials/condition), for both methods we used between 1 and 30 of the available trials per pseudotrial, and resampled each trial between 1 and 15 times. In Experiment 2 (45 trials/condition), we used between 1 and 15 of the available trials per pseudotrial, and resampled each trial between 1 and 10 times. Any trials that could not be used in a pseudotrial were left out of the classification. To ensure that the results were not dependent on the specific division of trials into pseudotrials, for each ‘subject’ and parameter set, we ran 100 iterations of the pseudotrial procedure and averaged the 100 resulting classification accuracies to give a single value for that subject and parameter set. Additional results across different numbers of iterations can be found in the supplementary material. For all simulations, we report the average decoding accuracy across subjects, the standard deviation, and their ratio (t-score against chance level decoding, 50%).

Results

Influence of trial averaging

In keeping with previous research, we show that averaging even a few trials together can be beneficial for classification performance. Importantly, for the dataset with no simulated difference between conditions (Figure 1, top panel) average decoding accuracy remained at 50% with the creation of pseudotrials. This reassures us that averaging the data cannot create effects where they are not present. For the data with simulated class distances greater than zero, we observed a complex pattern of results that varied with the number of trials per pseudotrial, classifier type, and class distance. Creating averages of even a few trials was helpful in most cases, resulting in increased decoding accuracy and higher t-statistics, despite the concurrent increase in standard deviation over subjects. However, including too many of the available trials in an average (reducing the number of datapoints available to the classifier) could be detrimental, as the increase in standard deviation outstripped the increase in average decoding,

resulting in reduced t-values. Using more than half of the available trials per pseudotrial was particularly detrimental for the Naïve Bayes classifier, which benefited from less averaging and more samples for the classifier. This was also the case for the SVM classifier, but to a lesser extent, and most visible when the class distance was the greatest.

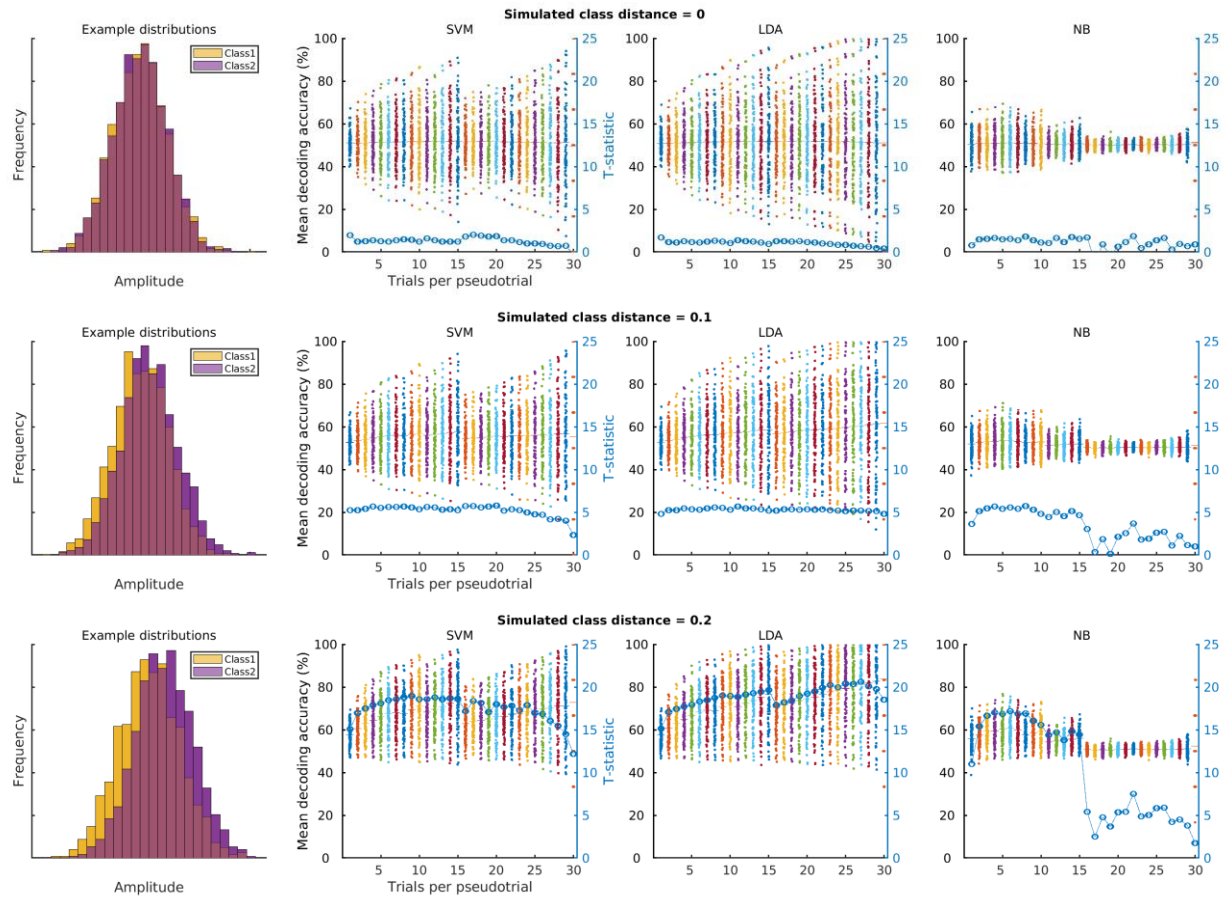


Figure 1. The influence of trial averaging, class distance, and classifier type on decoding accuracy and t-statistics. Rows correspond to the three simulated class distances 0, 0.1, and 0.2. In the left-hand column we plot example random normal distributions (randn) with added constants to illustrate the simulated class distances. This highlights that a class distance of 0.2 had the largest difference between the distribution means. Other columns correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). We simulated data from 100 subjects with 90 trials per condition. Pseudotrials were created separately within 3 ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated to pseudotrials across 100 iterations of pseudotrial creation for each subject. Each datapoint represents the average value for one subject.

Influence of trial resampling

Next, we examined the influence of trial resampling and its interaction with the number of trials per pseudotrial. Once again, with no simulated class differences, decoding remained at 50%, reassuring us that averaging and resampling the data cannot create effects where there are none (supplementary Figure 1a). For the data with a small simulated difference between conditions (distance = 0.1, Figure 2), we found that using around a third of the original trials per chunk to create each pseudotrial was sufficient to aid decoding (i.e., 10 of the 30 trials per chunk, or 11% of the total 90 trials). Combining this with a small amount of resampling further increased decoding accuracy, presumably because resampling allows more pseudotrials to be created. For example, with 30 trials per condition in each of the 3 chunks, averaging together 10 trials creates 3 pseudotrials per chunk if no resampling is used. With a resampling value of 2, each trial is included in 2 pseudotrials, meaning that 6 pseudotrials are created per chunk. While a resampling of 2 increased t-values for the Naïve Bayes classifier, this was not the case for SVM and LDA, where resampling mostly acted to stabilise the influence of the number of trials per pseudotrial. The results for a class distance of 0.2 were similar and can be found in supplementary Figure 1b.

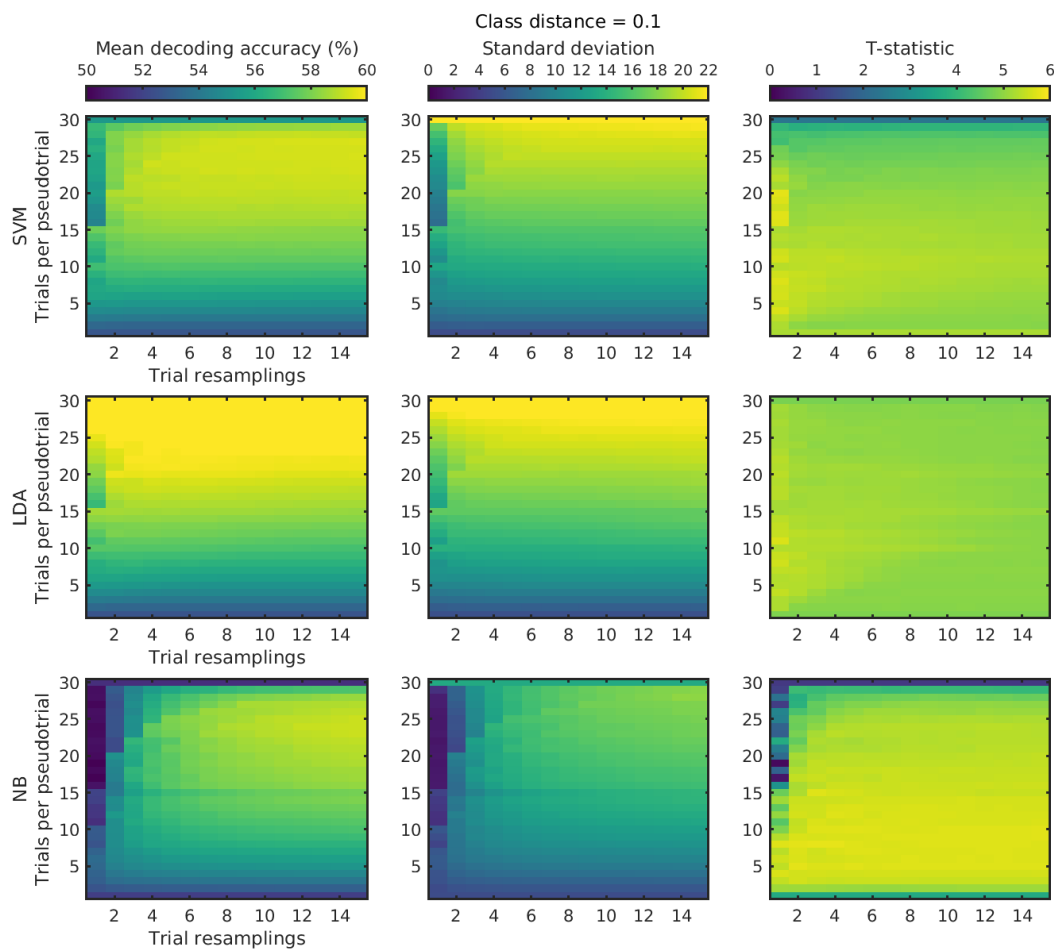


Figure 2. The influence of both averaging and resampling, on decoding accuracy, standard deviation, and t-statistics. Rows correspond to results from the three classifiers tested (SVM = support vector

machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0.1 (see supplementary Figures 1a and 1b for class distances of 0 and 0.2).

Influence of fewer trials per condition

Next, we checked whether the same principles would apply to an experiment with fewer trials per condition. Figure 3 demonstrates the effect of trial resampling on a smaller dataset with only 45 trials per condition, meaning 15 trials in each of the 3 chunks. Once again, a small amount of resampling combined trial averaging aided the classification particularly for the Naïve Bayes classifier. High values on either parameter was detrimental for classifier performance, and using up to a third of the original trials per chunk (i.e., 5 trials or less) per pseudotrial with a resampling of 2 was optimal.

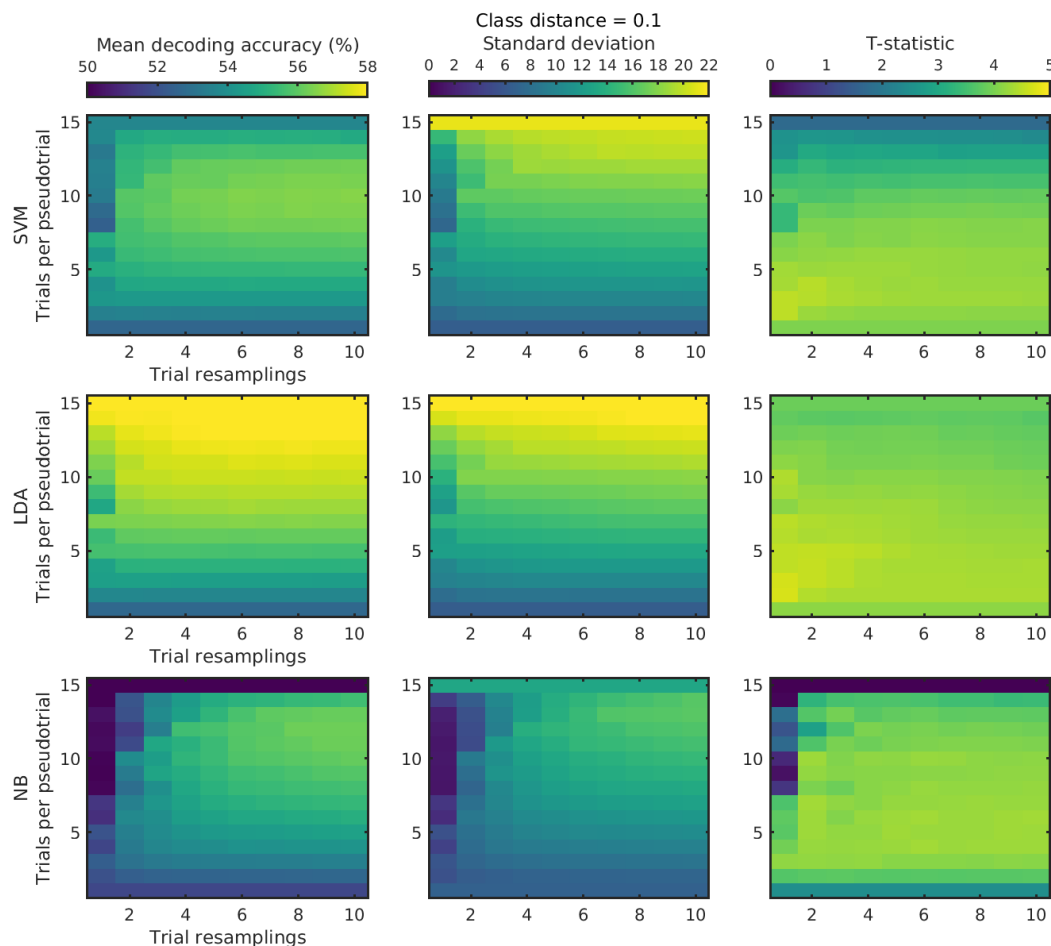


Figure 3. The influence of both averaging and resampling with fewer trials per condition. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects

with 45 trials per condition and a class distance of 0.1 (see supplementary Figure 2 for a class distance of 0.2).

Influence of a ‘one-pseudotrial-out’ decoding approach

The decoding analysis presented in Figures 1-3 utilised a chunking cross-validation procedure, in which trials were randomly allocated to one of three blocks before pseudotrials were created. Every analysis used a 3-fold cross validation, and only the number of trials available per block varied across averaging and resampling values. Next, we examined the alternative ‘one-pseudotrial-out’ procedure, where the number of folds was determined by the number of pseudotrials created. As shown in Figure 4, the maximum decoding accuracy and t-values achieved using the ‘one-pseudotrial-out’ method were slightly higher overall than the ‘3-block’ procedure, but had more variation across the parameter space. Once the number of folds reached the minimum of 2, the influence of the parameters was reduced. Less averaging was necessary to aid classification in the ‘one-pseudotrial-out’ approach, and the largest t-statistics were found across all classifiers when using roughly 5% of the total 90 trials per pseudotrial (i.e., 4 or 5 trials) with a resampling of 2.

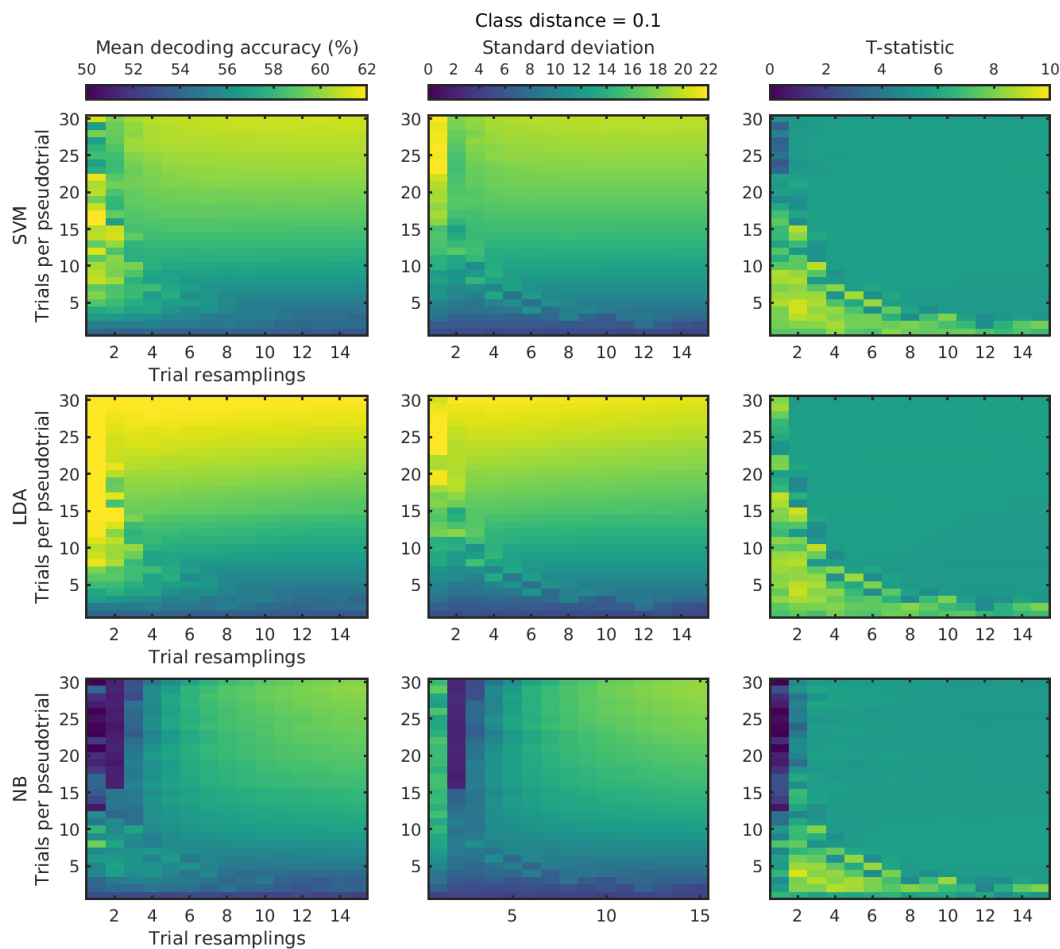


Figure 4. The influence of averaging and resampling using a ‘one-pseudotrial-out’ decoding approach. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA =

linear discriminant analysis, NB = Naïve Bayes). Here the number of cross-validation steps was determined by the number of trials per pseudotrial (the more trial averaging and resampling, the fewer possible steps), with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0.1 (see supplementary Figure 3 for a class distance of 0.2).

Discussion

We examined the influence of averaging and resampling across several parameters, including classifier type, class distance, number of trials available per condition, and the cross-validation procedure. When using a ‘3-chunk’ or ‘10-chunk’ approach, we found that using around a third of the original number of trials per chunk (or ~10% of all trials for 3-chunk, ~3% of all trials for 10-chunk) was sufficient to aid decoding. For the 3-chunk procedure, there was an additional stabilising effect provided by a low resampling value of 2. This was equivalent to creating 6 pseudotrials within each of the 3 chunks (18 pseudotrials in total), and was consistent for both experiment sizes. However, too much averaging could be detrimental, particularly for the Naïve Bayes classifier, and little was gained by using high resampling values. While increasing the number of trials used per pseudotrial generally increased decoding accuracy, it also increased the between-subject variance and therefore influenced the statistical outcome.

For the ‘one-pseudotrial-out’ cross-validation approach, fewer trials were needed per pseudotrial. Using around a sixth of the original number of trials per chunk (or ~5% of all trials) produced the largest t-statistics across all classifiers, when combined with a low resampling value of 2. Higher decoding accuracies and t-statistics were achieved for some parameters using the ‘one-pseudotrial-out’ approach compared to the ‘3-chunk’ version, but the different number of decoding steps resulted in more variation across the parameter space.

Although there were similarities across the three classifiers used, they responded differently across parameters, presumably due to the differences in their functions. Linear SVM separates classes by positioning a decision hyperplane in pattern space (Misaki et al., 2010). This hyperplane is chosen by maximising the distance to the patterns on either side, using the most informative data points that lie closest to the decision boundary (support vectors). Because of this, the SVM is not as influenced by changes in data points sitting away from the decision boundary. Therefore, SVM can perform well with limited data and will benefit most from having a few stable estimates near the decision boundary (Mur et al., 2009).

In LDA, the pattern space is constructed by maximising the between-class variance while minimising the within-class variance. The hyperplane is positioned in the middle of the class means, assuming that the two classes have Gaussian distributions and equal covariance. Therefore, a change in any data point will shift the decision boundary and potentially influence the classification result. Gaussian Naïve Bayes is similar to LDA, but also assumes that there are no correlations between pairs of data points within the same class (zero off-diagonal covariance). Having a low number of data points in each class distribution may more negatively impact the performance of Gaussian Naïve Bayes (Misaki et al., 2010), as demonstrated here.

Here, we focused on t-scores derived from decoding accuracy as a measure of classifier success, as this is the most common approach in MVPA studies. Other metrics of measuring

class separation have been proposed, such as cross-validated Mahalanobis distance (Walther et al., 2016). We believe that our results would generalise to such other measures of accuracy, as our results show that averaging increases separation between datapoints, evidenced by increased accuracy, which would similarly increase Mahalanobis distance. However, a full exploration of the effect of averaging on these other measures is outside the scope of the current paper.

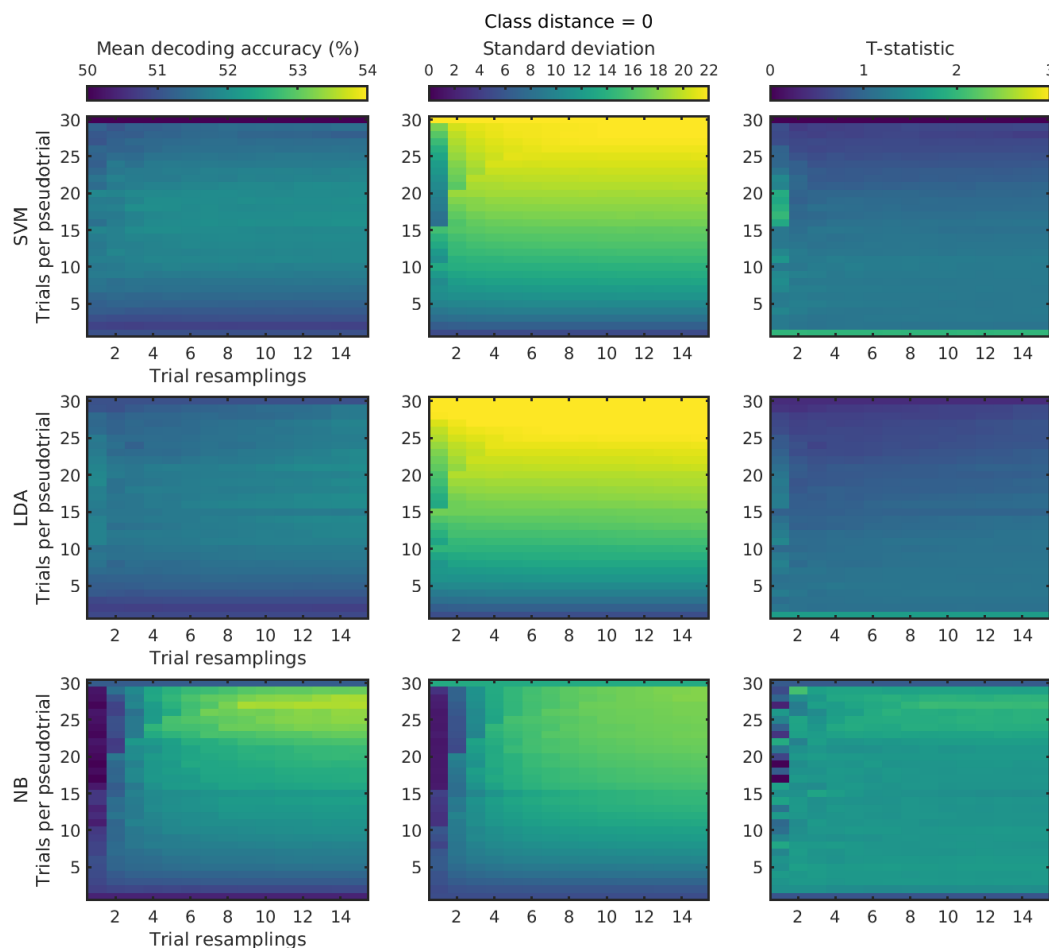
In summary, we found that modest trial averaging can improve decoding accuracy and associated t-statistics, and that a small amount of resampling helps to stabilise the benefit of doing so. However, only a low resampling value is helpful, and is not always necessary. In addition, the use of pseudotrials did not increase decoding accuracies when no effect was present. Although we provide general guidelines, the optimal parameter choice (particularly, the number of trials per pseudotrial) will be data and design specific, so we provide analysis code for others to run simulations based on their own design and hypothesised effects. We hope that our results and code can be used to inform future multivariate brain decoding studies.

Supplementary Material

Here we present additional results for class distances that were not included in the main text. This includes the influence of trial resampling, fewer trials per condition, and a one-pseudotrial-out approach. We also present additional analyses that were not addressed in the main text. This includes the influence of a smaller number of subjects, a smaller number of features, increasing the number of cross-validation steps to 10, and using different iterations of random trial allocation to pseudotrials.

Influence of trial resampling (additional class distances of 0 and 0.2)

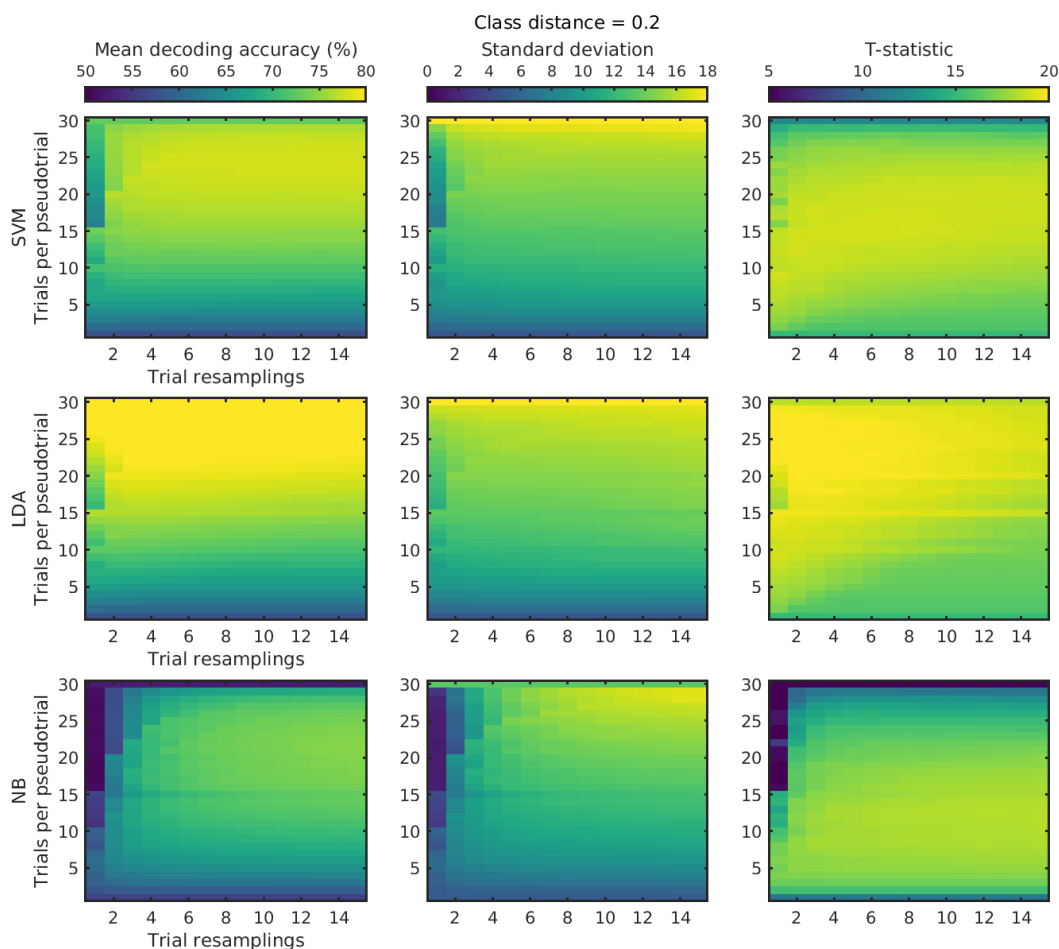
Here we examined the influence of trial resampling on decoding accuracy when there was no simulated class difference. This differs from main Figure 2 which displayed the results for a class distance of 0.1. Decoding performed on data with no simulated class differences remained at 50%, reassuring us that averaging and resampling the data cannot create effects where there is none (supplementary Figure 1a).



Supplementary Figure 1a. The influence of both averaging and resampling, on decoding accuracy, standard deviation, and t-statistics. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were

created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0 (see main Figure 2 for a class distance of 0.1).

For the data with a simulated difference between conditions of 0.1 (main Figure 2) we found that using around a third of the original trials per chunk to create each pseudotrial (i.e., 10 trials) was sufficient to aid decoding. Combining this with a small amount of resampling further increased decoding accuracy. As shown in supplementary Figure 1b, this was also true for the data a class distance of 0.2, when using the Naïve Bayes classifier. These parameters also performed well for the SVM and LDA classifiers, but SVM performed equally well with up to 50% of the original trials per pseudotrial (i.e., 15 trials) and a resampling value of 2. The LDA classifier performed well across a large range of parameters, providing that the resampling value was not higher than the number of trials per pseudotrial. With the larger simulated effect, the classifiers appeared to be more robust to the choice of averaging and resampling parameters.

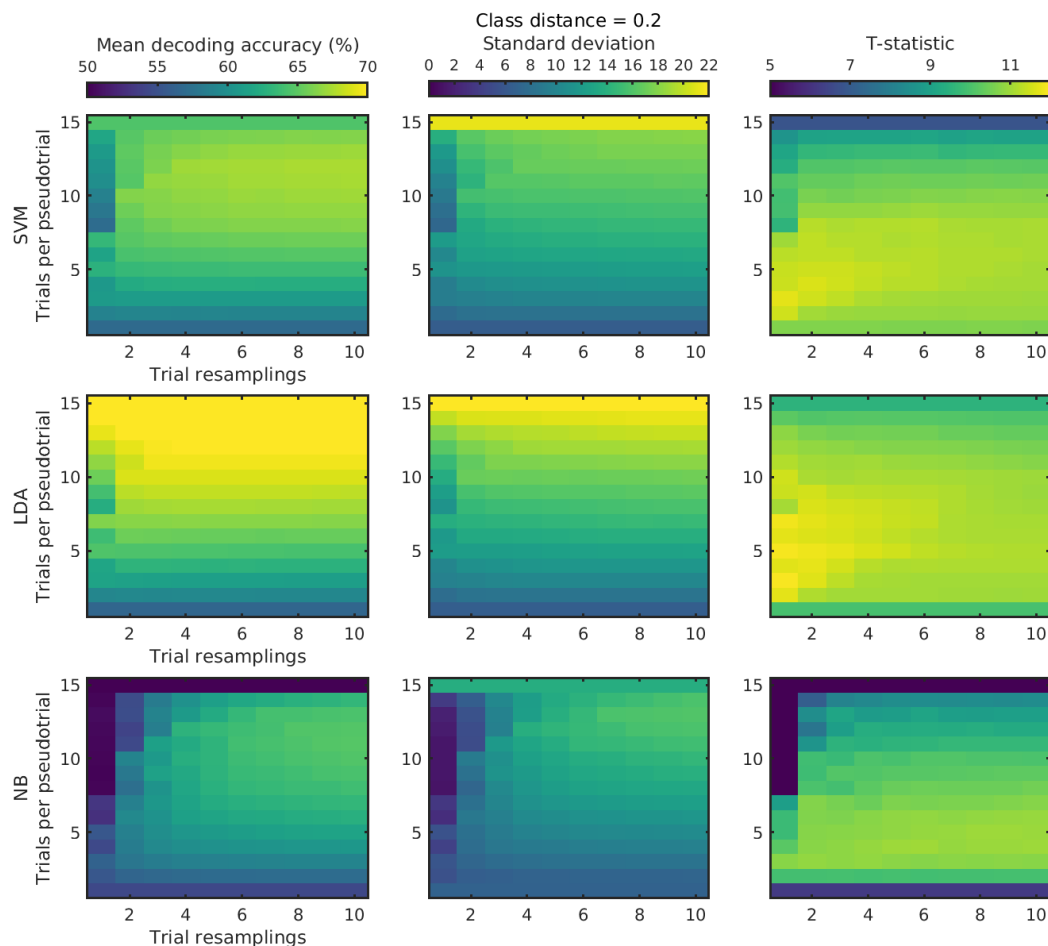


Supplementary Figure 1b. The influence of both averaging and resampling, on decoding accuracy, standard deviation, and t-statistics. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis). Pseudotrials were created separately

within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0.2 (see main Figure 2 for a class distance of 0.1).

Influence of fewer trials per condition (additional class distance of 0.2)

Here we examined the effect of trial resampling on a smaller dataset with only 45 trials per condition, meaning 15 trials in each of the 3 chunks. This differs from main Figure 3 which displayed the results for a class distance of 0.1. Once again, a small amount of resampling combined with trial averaging aided the classification. High values on either parameter was detrimental for classifier performance, and using up to a third of the original trials (i.e., 5 trials or less) per pseudotrial with a resampling of 2 was optimal for both a class distance of 0.1 (main Figure 3), and a class distance of 0.2 (supplementary Figure 2).

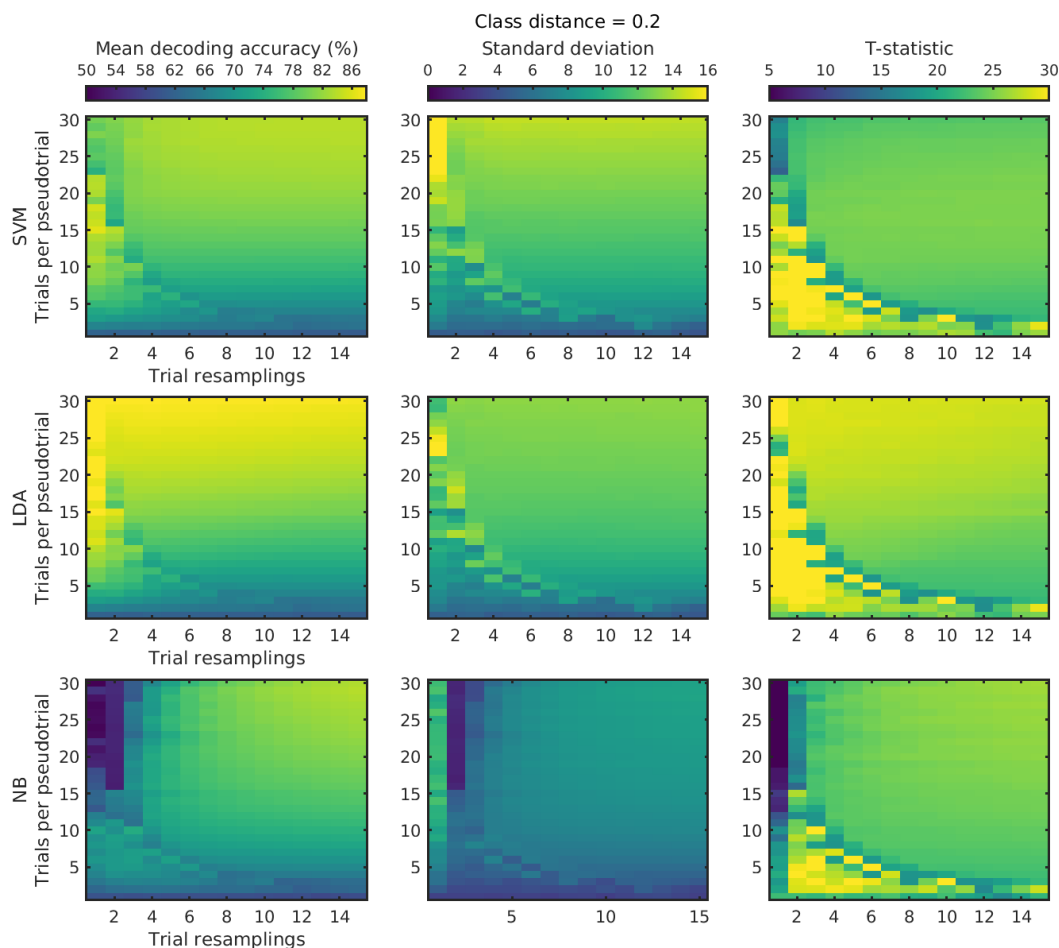


Supplementary Figure 2. The influence of both averaging and resampling with fewer trials per condition. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100

subjects with 45 trials per condition and a class distance of 0.2 (see main Figure 3 for a class distance of 0.1).

Influence of a ‘one-pseudotrial-out’ decoding approach (additional class distance of 0.2)

Here we examined a ‘one-pseudotrial-out’ procedure, where the number of folds was determined by the number of pseudotrials created, for a higher simulated class distance. This differs from main Figure 4 which displayed the results for a class distance of 0.1. For the Naïve Bayes classifier, high t-statistics were achieved when using roughly 15% of the original trials per pseudotrial (i.e., 4 or 5 trials) with a resampling of 2 (supplementary Figure 3). This is similar to the results using a class distance of 0.1 (main Figure 4). With a class distance of 0.2, the SVM classifier performed well with up to a third of the original trials (i.e., 10 trials) with a resampling of 2. For the LDA classifier, this could increase to half of the original trials (i.e., 15 trials) with a resampling of 2.

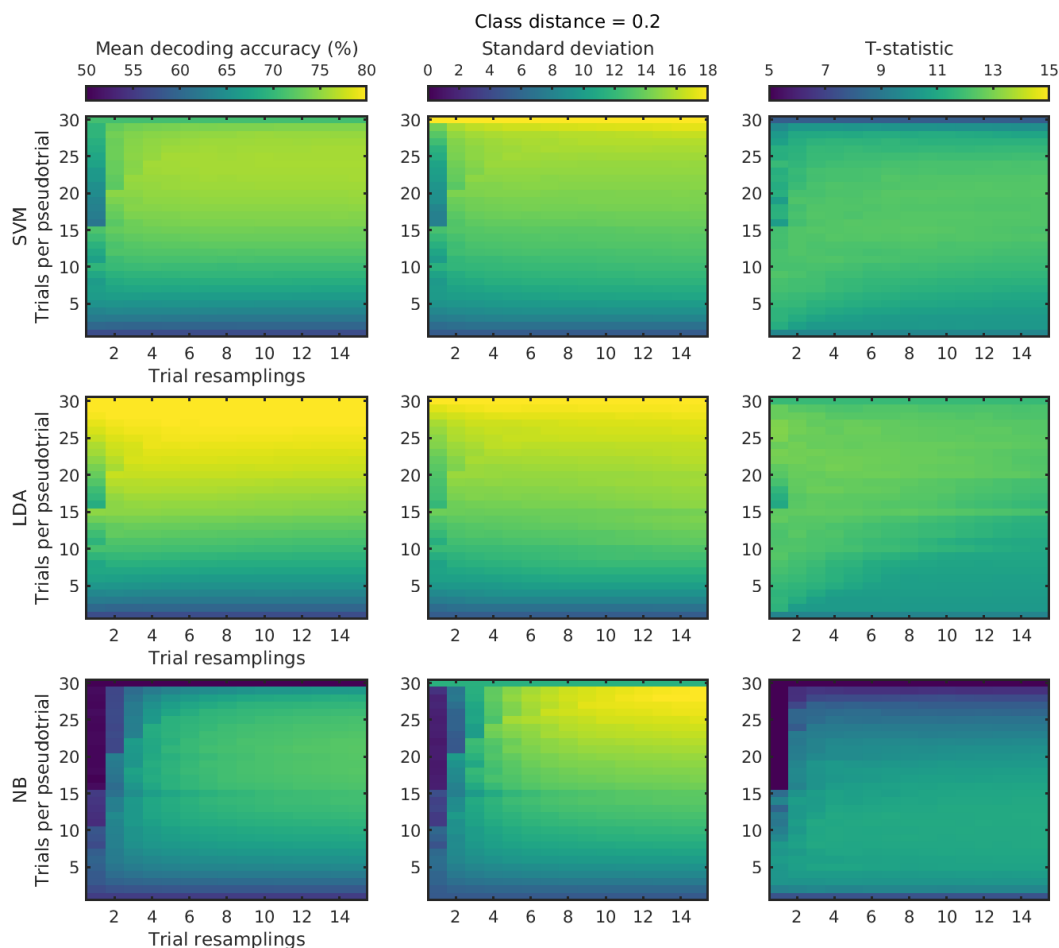


Supplementary Figure 3. The influence of averaging and resampling using a ‘one-pseudotrial-out’ decoding approach. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Here the number of cross-validation steps was determined by the number of trials per pseudotrial (the more trial averaging and resampling, the fewer possible steps), with trials randomly allocated across 100 iterations of pseudotrial creation.

For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0.2 (see main Figure 4 for a class distance of 0.1).

Influence of fewer ‘subjects’

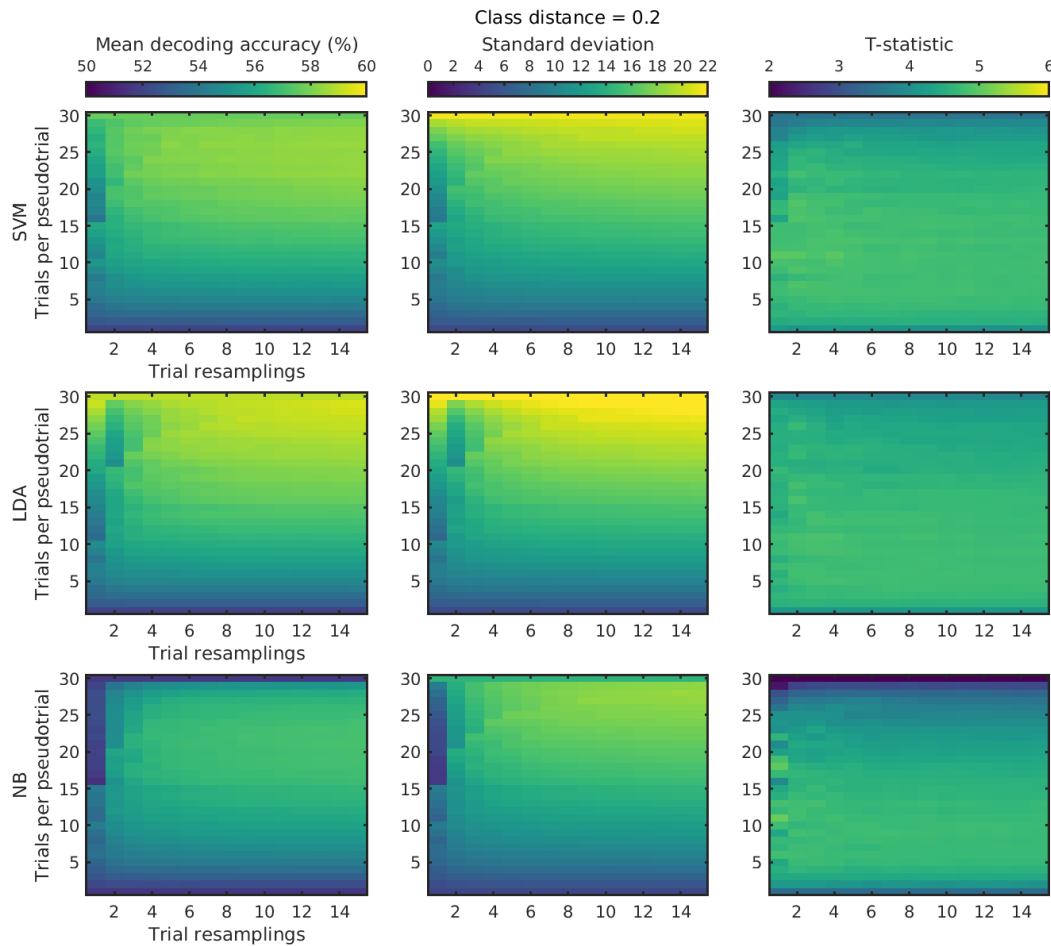
Here we examined the influence of reducing the number of simulated ‘subjects’ from 100 to 50, which was not examined within the main text. A similar pattern is found to the data in supplementary Figure 1b with 100 subjects a class distance of 0.2, although the overall performance is reduced. For the Naïve Bayes classifier, we found that using around a third of the original trials per chunk to create each pseudotrial (i.e., 10 trials) was sufficient to aid decoding. Combining this with a low resampling value of 2 further increased decoding accuracy. The SVM and LDA classifiers performed well with up to 50% of the original trials per pseudotrial (i.e., 15 trials) and a resampling value of 2.



Supplementary Figure 4. The influence of both averaging and resampling with fewer subjects. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 50 subjects with 90 trials per condition and a class distance of 0.2.

Influence of a smaller number of features

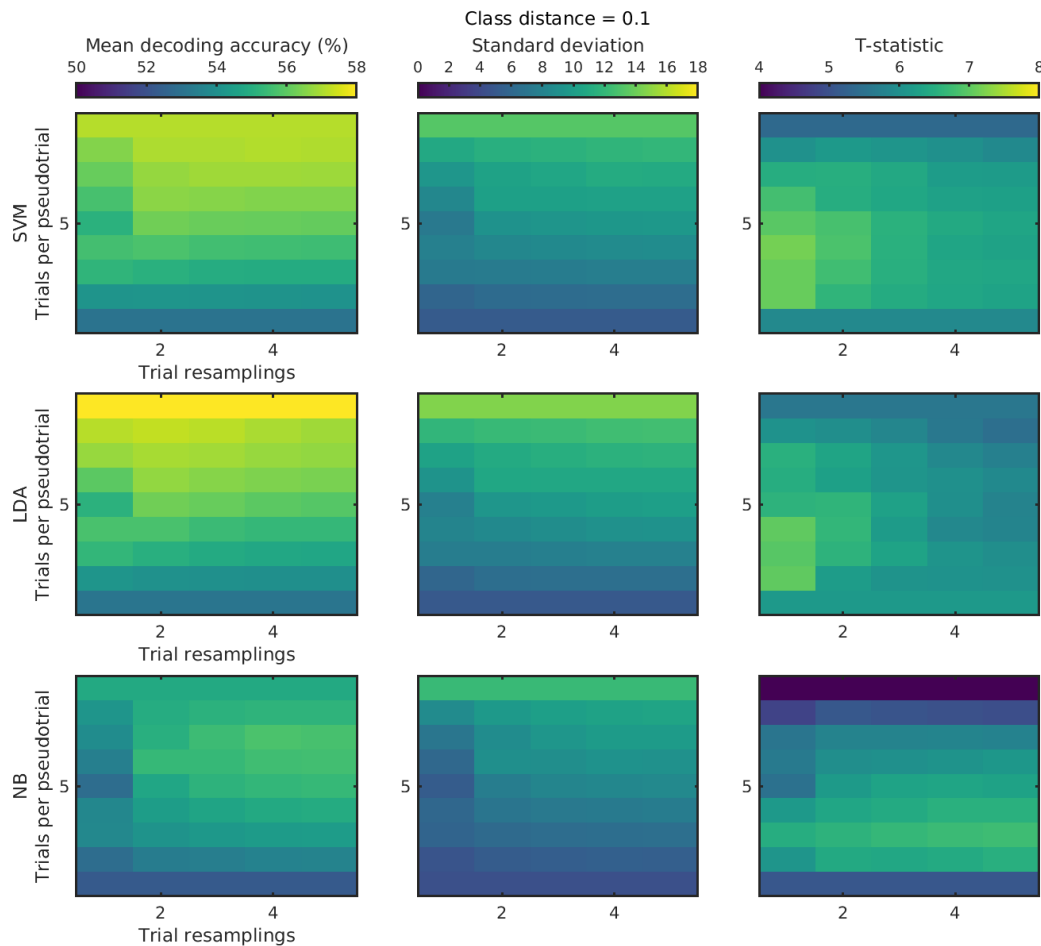
Here we examined the influence of reducing the number of simulated features from 700 to 6, which was not reported in the main text. The overall performance of the classifiers was reduced, as well as the influence of the parameters on t-statistics. However, using around a third of the original trials per pseudotrial and a resampling of 2 would still be a reasonable choice to optimise classification performance.



Supplementary Figure 5. The influence of both averaging and resampling with a smaller number of features (data size ‘small’ = 6 features). Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0.2.

Influence of increasing the number of cross-validation steps (10-fold)

Here we examined the influence of increasing the number of cross-validation steps from 3 to 10. As we simulated 90 trials per condition, there was a maximum of 9 trials per pseudotrial. As in the 3-chunk version (main Figure 2), using around a third of the original trials per chunk to create each pseudotrial (i.e., 3 trials) was sufficient to aid decoding. However, even for the Naïve Bayes classifier there was little to no benefit of resampling.

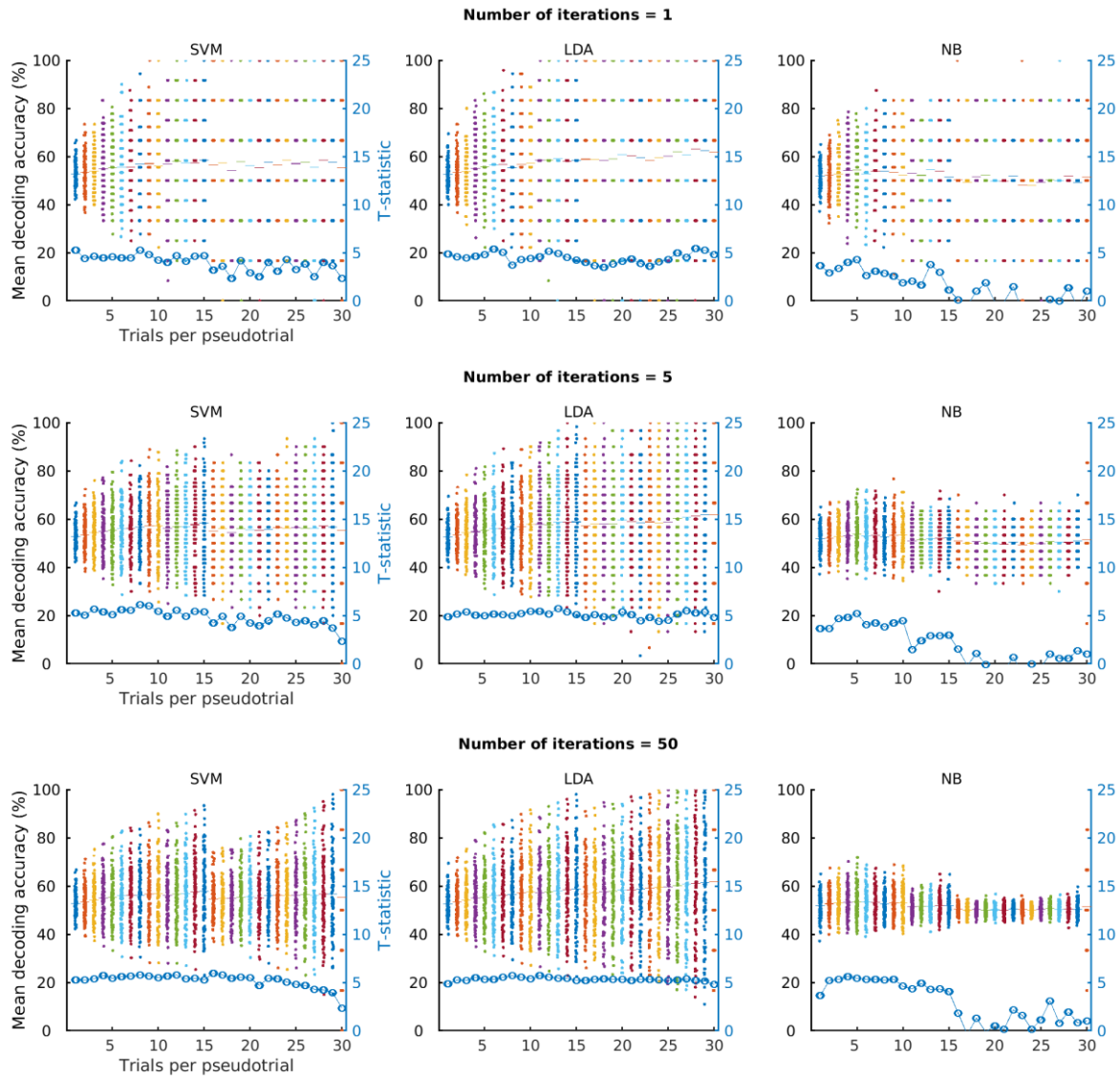


Supplementary Figure 6. The influence of averaging and resampling using a higher number of cross-validation steps. Rows correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 10 allocated ‘blocks’ of 9 trials each, facilitating a 10-fold cross-validation approach, with trials randomly allocated across 100 iterations of pseudotrial creation. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of 0.1

Influence of fewer iterations of random trial allocation

For the results reported in the main text, we ran 100 iterations of the pseudotrial procedure and averaged the 100 resulting classification accuracies. This was to ensure that the results were not dependent on the specific division of trails into pseudotrials, for each ‘subject’

and parameter set. Supplementary Figure 7 demonstrates the pattern of results achieved with one, five, and 50 iterations or random trial allocation. Without this iteration procedure, trial averaging quickly increases the between-subject variance, which is particularly detrimental for the performance of the Naïve Bayes classifier.



Supplementary Figure 7. The influence of fewer iterations of random trial allocation. Rows correspond to the number of iterations of random trial allocation that was used to create pseudotrials. Columns correspond to results from the three classifiers tested (SVM = support vector machine, LDA = linear discriminant analysis, NB = Naïve Bayes). Pseudotrials were created separately within 3 allocated ‘blocks’ of trials, facilitating a 3-fold cross-validation approach, with trials randomly allocated. For the results plotted here, we simulated data from 100 subjects with 90 trials per condition and a class distance of and 0.1.

Acknowledgements

This work was funded by the MRC Intramural funding, SUAG/093 G116768, and ARC Discovery Project, DP170101840. For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising from this submission. Thank you to Dr Arran Reader for their comments on an earlier version of the manuscript.

Data and Code Availability

Analysis scripts can be found at <https://osf.io/hjf75/>.

Author Contributions

CS: methodology, software, formal analysis, writing – original draft, visualisation. **TG:** conceptualisation, methodology, software, writing – editing and reviewing. **AW:** conceptualisation, methodology, software, writing – editing and reviewing, supervision, funding acquisition.

References

- Adam, K. C. S., Vogel, E. K., & Awh, E. (2020). Multivariate analysis reveals a generalizable human electrophysiological signature of working memory load. *Psychophysiology*, *57*(12), e13691. <https://doi.org/10.1111/psyp.13691>
- Bae, G.-Y., & Luck, S. J. (2018). Dissociable Decoding of Spatial Attention and Working Memory from EEG Oscillations and Sustained Potentials. *The Journal of Neuroscience*, *38*(2), 409–422. <https://doi.org/10.1523/JNEUROSCI.2860-17.2017>
- Bishop, C., & Nasrabadi, N. (2006). *Pattern Recognition and Machine Learning* / *SpringerLink* (Vol. 4). Springer. <https://link.springer.com/book/9780387310732>
- Duncan, D. H., van Moorselaar, D., & Theeuwes, J. (2023). Pinging the brain to reveal the hidden attentional priority map using encephalography. *Nature Communications*, *14*(1), Article 1. <https://doi.org/10.1038/s41467-023-40405-8>
- Foster, J. J., Bsaies, E. M., Jaffe, R. J., & Awh, E. (2017). Alpha-Band Activity Reveals Spontaneous Representations of Spatial Position in Visual Working Memory. *Current Biology*, *27*(20), 3216–3223.e6. <https://doi.org/10.1016/j.cub.2017.09.031>
- Goddard, E., Carlson, T. A., & Woolgar, A. (2022). Spatial and Feature-selective Attention Have Distinct, Interacting Effects on Population-level Tuning. *Journal of Cognitive Neuroscience*, *34*(2), 290–312. https://doi.org/10.1162/jocn_a_01796
- Grootswagers, T., Wardle, S. G., & Carlson, T. A. (2017). Decoding Dynamic Brain Patterns from Evoked Responses: A Tutorial on Multivariate Pattern Analysis Applied to Time Series Neuroimaging Data. *Journal of Cognitive Neuroscience*, *29*(4), 677–697. https://doi.org/10.1162/jocn_a_01068

- Haynes, J.-D., & Rees, G. (2006). Decoding mental states from brain activity in humans. *Nature Reviews Neuroscience*, 7(7), Article 7. <https://doi.org/10.1038/nrn1931>
- Hebart, M. N., & Baker, C. I. (2018). Deconstructing multivariate decoding for the study of brain function. *NeuroImage*, 180, 4–18. <https://doi.org/10.1016/j.neuroimage.2017.08.005>
- Hebart, M. N., Bankson, B. B., Harel, A., Baker, C. I., & Cichy, R. M. (2018). The representational dynamics of task and object processing in humans. *ELife*, 7, e32816. <https://doi.org/10.7554/eLife.32816>
- Isik, L., Meyers, E. M., Leibo, J. Z., & Poggio, T. (2014). The dynamics of invariant object recognition in the human visual system. *Journal of Neurophysiology*, 111(1), 91–102. <https://doi.org/10.1152/jn.00394.2013>
- Misaki, M., Kim, Y., Bandettini, P. A., & Kriegeskorte, N. (2010). Comparison of multivariate classifiers and response normalizations for pattern-information fMRI. *NeuroImage*, 53(1), 103–118. <https://doi.org/10.1016/j.neuroimage.2010.05.051>
- Mur, M., Bandettini, P. A., & Kriegeskorte, N. (2009). Revealing representational content with pattern-information fMRI—an introductory guide. *Social Cognitive and Affective Neuroscience*, 4(1), 101–109. <https://doi.org/10.1093/scan/nsn044>
- Oosterhof, N. N., Connolly, A. C., & Haxby, J. V. (2016). CoSMoMVPA: Multi-Modal Multivariate Pattern Analysis of Neuroimaging Data in Matlab/GNU Octave. *Frontiers in Neuroinformatics*, 10. <https://www.frontiersin.org/articles/10.3389/fninf.2016.00027>
- Pereira, F., Mitchell, T., & Botvinick, M. (2009). Machine learning classifiers and fMRI: A tutorial overview. *NeuroImage*, 45(1, Supplement 1), S199–S209. <https://doi.org/10.1016/j.neuroimage.2008.11.007>

Petit, S., Brown, A., Jessen, E. T., & Woolgar, A. (2023). How robustly do multivariate EEG patterns track individual-subject lexico-semantic processing of visual stimuli?

Language, Cognition and Neuroscience, 0(0), 1–15.

<https://doi.org/10.1080/23273798.2023.2177315>

Tuckute, G., Hansen, S. T., Pedersen, N., Steenstrup, D., & Hansen, L. K. (2019). Single-

Trial Decoding of Scalp EEG under Natural Conditions. *Computational Intelligence*

and Neuroscience, 2019, 1–11. <https://doi.org/10.1155/2019/9210785>

Valente, G., Castellanos, A. L., Hausfeld, L., De Martino, F., & Formisano, E. (2021). Cross-

validation and permutations in MVPA: Validity of permutation strategies and power

of cross-validation schemes. *NeuroImage*, 238, 118145.

<https://doi.org/10.1016/j.neuroimage.2021.118145>

Varoquaux, G., Raamana, P. R., Engemann, D. A., Hoyos-Idrobo, A., Schwartz, Y., &

Thirion, B. (2017). Assessing and tuning brain decoders: Cross-validation, caveats, and guidelines. *NeuroImage*, 145, 166–179.

<https://doi.org/10.1016/j.neuroimage.2016.10.038>

Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., & Diedrichsen, J. (2016).

Reliability of dissimilarity measures for multi-voxel pattern analysis. *NeuroImage*,

137, 188–200. <https://doi.org/10.1016/j.neuroimage.2015.12.012>