



Class imbalance should not throw you off balance: Choosing the right classifiers and performance metrics for brain decoding with imbalanced data

Philipp Thölke^{a,b,*}, Yorguin-Jose Mantilla-Ramos^{a,c}, Hamza Abdelhedi^a, Charlotte Maschke^{a,d}, Arthur Dehgan^{a,h}, Yann Harel^a, Anirudha Kementur^a, Loubna Mekki Berrada^a, Myriam Sahraoui^a, Tammy Young^{a,e}, Antoine Bellemare Pépin^{a,f}, Clara El Khantour^a, Mathieu Landry^a, Annalisa Pascarella^g, Vanessa Hadid^a, Etienne Combrisson^h, Jordan O'Byrne^a, Karim Jerbi^{a,i,j}

^a Cognitive and Computational Neuroscience Laboratory (CoCo Lab), University of Montreal, 2900, boul. Edouard-Montpetit, Montreal, H3T 1J4, Quebec, Canada

^b Institute of Cognitive Science, Osnabrück University, Neuer Graben 29/Schloss, Osnabrück, 49074, Lower Saxony, Germany

^c Neuropsychology and Behavior Group (GRUNECO), Faculty of Medicine, Universidad de Antioquia, 53-108, Medellín, Aranjuez, Medellín, 050010, Colombia

^d Integrated Program in Neuroscience, McGill University, 1033 Pine Ave, Montreal, H3A 0G4, Canada

^e Department of Computing Science, University of Alberta, 116 St & 85 Ave, Edmonton, T6G 2R3, AB, Canada

^f Department of Music, Concordia University, 1550 De Maisonneuve Blvd. W., Montreal, H3H 1G8, QC, Canada

^g Institute for Applied Mathematics Mauro Picone, National Research Council, Roma, Italy

^h Institut de Neurosciences de la Timone (INT), CNRS, Aix Marseille University, Marseille, 13005, France

ⁱ Mila (Quebec Machine Learning Institute), 6666 Rue Saint-Urbain, Montreal, H2S 3H1, QC, Canada

^j UNIQUE Centre (Quebec Neuro-AI Research Centre), 3744 rue Jean-Brillant, Montreal, H3T 1P1, QC, Canada

ARTICLE INFO

Keywords:

Class imbalance
Machine learning
Classification
Performance metrics
Electroencephalography
Magnetoencephalography
Brain decoding
Balanced accuracy

ABSTRACT

Machine learning (ML) is increasingly used in cognitive, computational and clinical neuroscience. The reliable and efficient application of ML requires a sound understanding of its subtleties and limitations. Training ML models on datasets with imbalanced classes is a particularly common problem, and it can have severe consequences if not adequately addressed. With the neuroscience ML user in mind, this paper provides a didactic assessment of the class imbalance problem and illustrates its impact through systematic manipulation of data imbalance ratios in (i) simulated data and (ii) brain data recorded with electroencephalography (EEG), magnetoencephalography (MEG) and functional magnetic resonance imaging (fMRI). Our results illustrate how the widely-used Accuracy (Acc) metric, which measures the overall proportion of successful predictions, yields misleadingly high performances, as class imbalance increases. Because Acc weights the per-class ratios of correct predictions proportionally to class size, it largely disregards the performance on the minority class. A binary classification model that learns to systematically vote for the majority class will yield an artificially high decoding accuracy that directly reflects the imbalance between the two classes, rather than any genuine generalizable ability to discriminate between them. We show that other evaluation metrics such as the Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC), and the less common Balanced Accuracy (BAcc) metric - defined as the arithmetic mean between sensitivity and specificity, provide more reliable performance evaluations for imbalanced data. Our findings also highlight the robustness of Random Forest (RF), and the benefits of using stratified cross-validation and hyperparameter optimization to tackle data imbalance. Critically, for neuroscience ML applications that seek to minimize overall classification error, we recommend the routine use of BAcc, which in the specific case of balanced data is equivalent to using standard Acc, and readily extends to multi-class settings. Importantly, we present a list of recommendations for dealing with imbalanced data, as well as open-source code to allow the neuroscience community to replicate and extend our observations and explore alternative approaches to coping with imbalanced data.

* Corresponding author.

E-mail address: philipp.thoelke@posteo.de (P. Thölke).

<https://doi.org/10.1016/j.neuroimage.2023.120253>.

Received 18 April 2023; Received in revised form 5 June 2023; Accepted 26 June 2023

Available online 28 June 2023.

1053-8119/© 2023 Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

1. Introduction

The rise of artificial intelligence (AI) in the last decade has led to important breakthroughs across many areas of science, including neuroscience and neuroimaging. New synergies between neuroscience and AI promise to drive both fields forward (Gershman et al., 2015; Hassabis et al., 2017; Helmstaedter, 2015; Macpherson et al., 2021; Richards et al., 2019). In particular, machine learning is increasingly used both to model and to classify brain data (Yang and Wang, 2020), with applications ranging from cognitive and systems neuroscience (Fong et al., 2018) to clinical brain imaging (Buchlak et al., 2021; Myszczyńska et al., 2020). As a result, machine learning is steadily turning into a fundamental tool for neuroscientists (Glaser et al., 2019). As is the case with all methodological frameworks, machine learning comes with a set of subtleties and pitfalls. Being aware of these limitations and knowing how to handle them properly can be challenging, especially in research domains where machine learning is not yet adequately and systematically covered during training. The issue of data imbalance (He and Garcia, 2009; Sun et al., 2009) is a perfect example of an important problem that is generally well understood in the field of data science, but not always properly appreciated and tackled in neuroscience and neuroimaging. This technical note provides (1) a didactic description of the pitfalls associated with using skewed datasets in supervised machine learning, (2) a detailed assessment of the impact of varying the degree of class imbalance on classifier models and their performance using synthetic and real data, (3) concrete recommendations for mitigating the adverse effects of imbalanced data, and (4) open-source code to replicate the present work and extend it to other methods and metrics.

In binary classification problems, data imbalance occurs whenever the number of observations from one class (majority class) is higher than the number of observations from the other class (minority class) (He and Garcia, 2009; Sun et al., 2009). This problem is commonly encountered in cognitive neuroscience and in clinical applications, where observations for the target class (e.g. patients with neurological disorders) are often much harder to come by than for the control class (e.g. cognitively healthy individuals), leading to datasets with many more control observations than target observations (Krawczyk, 2016; Sun et al., 2009). Additional care has to be taken when evaluating the performance of diagnostic tests on rare conditions (Varoquaux and Colliot, 2022).

What makes imbalanced data problematic? When faced with highly skewed data, a classifier can achieve a high decoding accuracy merely by systematically and blindly voting for the majority class (Sun et al., 2009). For example, if an image classifier is asked to discriminate pictures of crows versus ravens, but only one out of twenty images in the training and test sets are ravens and the rest are crows, then the algorithm can (and likely will) achieve 95% accuracy simply by calling everything it sees a crow—though no discrimination or classification can rightly said to have been accomplished. In other words, given the opportunity, an algorithm will tend to bypass more complex feature analysis simply by “playing the odds”, which is indistinguishable from actual classification when only focusing on e.g. Accuracy as a performance metric.

The most common approach to avoid this problem is to enforce balanced data. One way to do this is by undersampling, i.e. by removing observations from the majority class until a balance is reached (Chawla et al., 2004; Liu et al., 2008), and repeating the process through bootstrapping. However, this comes at the cost of reducing the sample size, increasing the signal-to-noise ratio, which can be detrimental to the classification. Alternatively, one can oversample the minority class by duplicating or interpolating observations (Chawla et al., 2002; Fernández et al., 2018; Graa and Rekik, 2019) (Fig. 1a), though this comes with a higher risk of overfitting and introducing noise (Tan et al., 2007).

It may also be possible to dispense with undersampling or oversampling, and the problems they create, and to cope with the imbalanced data. In this case however, a number of additional considerations are necessary to avoid spurious results (Haixiang et al., 2017). These include

the judicious choice of the type of classifier and the performance metric to be used. Additionally, when deploying a model validation scheme, special care must be taken to reflect the imbalance in the main data, such as by using Stratified K-Fold cross-validation (Fig. 1b). While these best practices are commonly applied in the machine learning community, they are not as widely adopted by the neuroscience and neuroimaging fields, likely due to the little information that exists, targeting neuroscientists, on how each of these different factors interact with data imbalance in a neuroscience context. Nevertheless, the importance of considering class imbalance has been highlighted in several brain decoding studies, and appropriate metrics have been made available in some toolboxes (Bode et al., 2019; Fahrenfort et al., 2018; Pereira et al., 2009). Recommended metrics include the Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC), d-prime (a metric related to ROC-AUC, commonly used in psychology) (Das and Geisler, 2021; Hebart et al., 2015), balanced accuracy (Grootswagers et al., 2017) and tuned loss functions (Lemm et al., 2011). Yet, a dedicated account and systematic quantification of the effect of data imbalance for the neuroscience community is still missing. In particular, it is useful to have a didactic assessment of the impact of data imbalance (using both synthetic and real data) across varying degrees of imbalance, types of classifiers, hyperparameter choices, training, cross-validation and significance testing schemes. In this paper, we aim to provide a straightforward and practical demonstration of this multifaceted problem by using simulated data, as well as real-world electrophysiological and fMRI recordings. More specifically, we examine the behavior of four prominent metrics (Accuracy (Acc), Area Under the ROC Curve (AUC), Balanced Accuracy (BAcc) (Brodersen et al., 2010a; Kelleher et al., 2015), and F1; Table 1) across four widely-used classifiers (Logistic Regression (LR) (Cox, 1958), Linear Discriminant Analysis (LDA) (Fukunaga, 1993), Support Vector Machine (SVM) (Cortes and Vapnik, 1995), and Random Forest (RF) (Breiman, 2001)), as we gradually increment data imbalance.

The topic of data imbalance, also often referred to as class or domain imbalance, has been addressed in previous work and online resources (Haixiang et al., 2017; Sun et al., 2009), primarily within the computer science community. Here, we tailor our examples, explanations and recommendations, as well as our open-source code, to the neuroscience researcher or trainee with an interest in applying machine learning to neuroimaging data.

2. Methods and materials

To explore the effect of data imbalance on different classification algorithms and performance metrics, we developed a custom open-source analysis pipeline, which systematically manipulates class imbalance (Fig. 1c). We herein first describe the analysis pipeline and secondary analysis, and then describe the five datasets used in this study (i.e. three types of simulated data, one EEG dataset, and one MEG dataset).

2.1. Synthetic data

To evaluate the impact of class imbalance in a controlled environment we generated synthetic data, consisting of 1000 random samples (at perfect balance) from two Gaussian distributions. We explored three different scenarios by modifying the amount of overlap between the two distributions, i.e. changing the distance between the means μ_1 and μ_2 of the two distributions, while keeping the standard deviation σ_1 and σ_2 constant at 1. In the first scenario, both classes came from the same distributions ($|\mu_0 - \mu_1| = 0$; Fig. 2a) and are therefore impossible to classify. In the second scenario, the two distributions were mostly overlapping ($|\mu_0 - \mu_1| = 1$; Fig. 2f), simulating a hard classification task. In the third scenario, the two distributions had a minimal overlap ($|\mu_0 - \mu_1| = 3$; Fig. 2k), which illustrates an easy classification task.

Note that this dataset serves a didactic purpose only, as we limit the simulated data to a single feature. While oftentimes one has to deal with

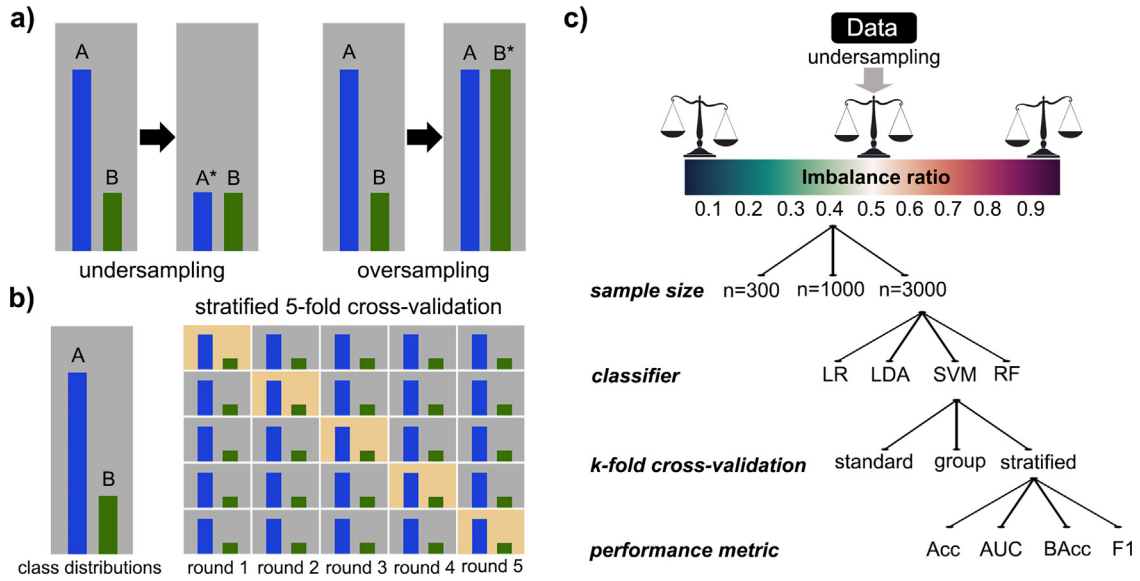


Fig. 1. **a)** Methods to balance imbalanced data in order to avoid biases in machine learning. In undersampling (left), a subset of the overrepresented group (dataset A) is chosen. In oversampling (right), samples from the underrepresented group (dataset B) are duplicated or artificially augmented. **b)** Illustration of Stratified K-Fold cross validation (K=5). Instead of randomly choosing subsamples for every fold, this technique maintains the balance of the original data over all folds. This technique helps reduce biases and large variance in cross-validation. **c)** Illustration of the overall analysis framework of experiments performed in this paper. Various degrees of class imbalance were manually generated by undersampling the data. For a set of sample sizes, we performed binary classification using four widely used algorithms, three K-fold cross-validation methods, and four evaluation metrics (Acc, AUC, BAcc, and F1).

multiple features, this example illustrates the problem at hand in a simplistic way and makes the concepts easy to grasp. We further extend our analysis to real-world examples using multiple features (Section 2.4).

2.2. Primary analysis pipeline

The analysis pipeline was developed specifically for binary classification problems. Its primary purpose is to generate scores for different metrics across a range of imbalance ratios, using a list of classifiers and cross-validation schemes (Fig. 1c).

In order to estimate the chance level of correct classification, given the configuration of dataset and performance metric, the pipeline performs permutation tests (Ojala and Garriga, 2010) (repeatedly training and evaluating a classifier on the same dataset but with randomly permuted labels). Generating data-driven chance level is necessary because the theoretical binary classification chance level of 0.5 could be incorrect when performing binary classification on a dataset which has imbalanced classes or a small sample size. Furthermore, different performance metrics can lead to different estimates of the chance level.

In addition to that, we repeated the experiments 10 times with different random seeds to estimate the degree of variance across cross-validation splits. We generally report the mean performance across the 10 repetitions and indicate the standard deviation as a shaded area around the mean.

To assess the impact of the metric with which classifiers are evaluated, we explored a range of classification metrics. These include Accuracy (Acc) (Luque et al., 2019), Balanced Accuracy (BAcc) (Brodersen et al., 2010b), Area Under the ROC Curve (AUC) (Gong, 2021), and F1 (Wang et al., 2017). We made sure to include the most frequently used metrics, as well as variations specifically designed to tackle evaluation of prediction on imbalanced data. See Table 1 for an overview of classification metrics. It is interesting to note here that the Area under the Curve (AUC) of the Receiver Operating Characteristic (ROC) measures the integral of true positive rate against the false positive rate across all decision thresholds. This means that AUC is invariant to modifying the decision threshold, which is sometimes used to combat the influence of class imbalance.

In terms of classifiers, we included Logistic Regression (LR) (Cox, 1958), Linear Discriminant Analysis (LDA) (Fukunaga, 1993), Sup-

Table 1

Overview of evaluation metrics. True positives (TP): instances that are positives and are classified as positives. False positive (FP): instances that are negatives and are classified as positives. False negative (FN): instances that are positives and are classified as negatives. True negatives (TN): instances that are negatives and are classified as negatives. PPV, Positive Predictive Value. NPV, Negative Predictive Value.

Metric	Definition	Formula
Precision/ PPV	Correct positive predictions divided by all positive predictions.	$\frac{TP}{TP+FP}$
Recall/ Sensitivity	Correct positive predictions divided by all positive samples.	$\frac{TP}{TP+FN}$
Specificity	Correct negative predictions divided by all negative samples.	$\frac{TN}{TN+FP}$
NPV	Correct negative predictions divided by all negative predictions.	$\frac{TN}{TN+FN}$
F1 score	The harmonic mean of precision and recall.	$2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
Accuracy	Proportion of correct predictions among all samples.	$\frac{TP+TN}{TP+FP+TN+FN}$
Balanced Accuracy	Mean of recall and specificity, i.e. average per-class accuracy.	$0.5 \cdot \left(\frac{TP}{TP+FN} + \frac{TN}{TN+FP} \right)$
AUC	Area Under ROC Curve, which plots true positive against false positive rate for all decision thresholds.	$\int_0^1 \text{ROC}$

port Vector Machine (SVM) (Cortes and Vapnik, 1995), and Random Forest (RF) (Breiman, 2001). We chose these models because they are among the most widely used in the neuroscience community, and because they represent a variety of approaches. The SVM was used with a radial basis function (RBF) kernel, making it a powerful non-linear classification algorithm. In addition we also assessed the impact of class imbalance on the commonly used linear SVM classifier (see Supplementary Fig. 2). RF was of specific interest as it is a tree-based ensemble model expected to be better at handling class imbalance than the other methods.

We essentially used the default hyperparameters as defined in the scikit-learn library (Pedregosa et al., 2011), however, we reduced the number of Random Forest estimators from 100 to 25 to better suit the low number of features in our experiments. Moreover, computing AUC requires predicting a continuous classification score for each sample. While LR, LDA and RF provide probability predictions out of the box, SVM supports two ways of estimating these prediction scores: probability calibration through an internal cross-validation on the training set, or using the signed distance from the fitted hyperplane. We found no substantial difference between both techniques and used the probability calibration method throughout the analysis. Note that AUC and Balanced Accuracy would be equivalent if we used binary predictions instead of class probabilities.

Cross-validation was performed using the Stratified K-fold or Stratified Group K-fold strategy (5 folds), depending on the presence or absence of group/subject information in the data. In a typical data-driven neuroscience decoding task, group labels help separate data from different subjects and add a measure of generalisation performance to new subjects to the evaluation process.

To simulate different amounts of imbalance in the class distribution we artificially limited the number of samples for both classes separately. We used a range of imbalance ratios from 0.1 (9:1 balance between the two classes) to 0.9 (1:9 balance) with 25 linearly spaced intermediate ratios, which provided a good trade-off between speed and performance. For a dataset with 100 data points (50 in either class), for example, we ran experiments with the following class distributions: 50:5, 50:7, ..., 50:50, ..., 7:50, 5:50. Imbalance was achieved by undersampling (dropping samples) either one of the two classes. It is important to note here that the sample size used to fit the classifiers decreased with increasing levels of imbalance. This limitation comes in part from the limited amount of data in our EEG and MEG analysis. We investigated the effect of sample size on our analysis in later stages of analysis, ensuring that it does not interfere with our results.

2.3. Secondary analysis

Additionally, we explored the effect of data imbalance as a function of 1) the selection of hyperparameters, 2) the size of the dataset, 3) the type of cross-validation, 4) the effect of a balanced hold-out test set and 5) the impact of class imbalance on statistical significance testing. To simplify our approach, we only performed analysis 2–5 using SVM with an RBF kernel and evaluated it using Acc. We chose SVM specifically because we expected this algorithm to display important effects of class imbalance on performance.

2.3.1. Effect of hyperparameters

To assess the putative effect of hyperparameters, we explored those that are expected to have a significant impact on the robustness of classifiers with respect to imbalanced data (Zhu et al., 2018). We used synthetic data (Section 2.1) with 1000 samples (at perfect balance) and a distance of one between the two Gaussian distributions and evaluated the effect of the selected hyperparameters using Balanced Accuracy. This allows us to track improvements in robustness, which would manifest as a flatter curve of classification scores across imbalance ratios.

We limited this experiment to hyperparameters implemented in scikit-learn as this is one of the most commonly used libraries. Logis-

tic Regression, SVM, and Random Forest all implement an automatic class-weighting algorithm to deal with imbalanced data, which can be enabled by setting the *class_weight* parameter to *balanced*. This approach weights the influence of each sample according to the inverse frequency of the corresponding class, thereby decreasing the impact of the majority class. This class-weighting technique, also known as cost-sensitive learning, penalizes the model less for errors made on examples from the majority class and more for errors made on the minority class. The Random Forest additionally has a *balanced_subsample* option, which applies the same weighting on the level of individual trees instead of globally for the full model. In addition to class weighting, we explored changing the minimum size of leaf nodes as a fraction of all samples (*min_weight_fraction_leaf*). As the fractional size of the leaf nodes depends on class distribution a large enough value ensures that leaf nodes will be more representative of differences between classes instead of simply voting for the majority class. We explored values of 0.1 and 0.4 for this parameter, serving as examples for weak and strong regularization.

While these hyperparameter optimization strategies are readily available for LR, SVM and RF in their respective scikit-learn implementations, tackling class imbalance by changing model parameters in the case of LDA is less straightforward. One simple strategy would be to modify the intercept of the decision hyperplane according to the rate of class imbalance. While this can be achieved through the *priors* hyperparameter in scikit-learn, exploration of this hyperparameter in more detail goes beyond the scope of this article. More generally, while hyperparameter optimization is not done systematically in brain decoding work, it is likely to become more common as ML continues to be increasingly used in neuroscience. Hence, understanding its impact on the issue of class imbalance could become increasingly relevant.

2.3.2. Effect of sample size

To test the impact of dataset size on robustness to imbalance we evaluated classifier performance across imbalance ratios using synthetic data with sample sizes $N=300$, $N=1000$ and $N=3000$ before undersampling. The data of the two classes was sampled from two Gaussian distributions with a distance of one (Section 2.1). This experiment further allows us to shed light on the issue of decreasing sample size with increasing levels of class imbalance, which results from the technique we used to generate imbalanced datasets.

2.3.3. Effect of cross-validation

In order to assess the influence of the cross-validation scheme on different metrics when using imbalanced data, we tested K-Fold and Stratified K-Fold cross-validation on synthetic data. This difference will likely appear on smaller sample size since lower sample sizes increase the likeliness of having one class absent from a fold when using K-Fold without stratification.

To assess the impact of the choice of cross-validation approach, we trained an SVM classifier on the synthetic data with a distance of one between the means of both distributions (Section 2.1) and chose to only use 50 samples per class before unbalancing. This analysis was repeated with 40 different seeds in order to assess the robustness of the effects we hypothesise.

2.3.4. Effect of balanced hold-out set

While so far all cross-validation splits came from the training data distribution, thereby replicating class imbalance, we also decided to explore the effects of training on an imbalanced dataset and performing validation on a balanced subset. Here, we trained an SVM on 1000 samples (at perfect balance) of synthetic data (Section 2.1) with a distance of one between the means of both classes. The balanced hold-out set was created by taking a 10% split of the full dataset before artificially generating an imbalanced training set. This analysis aims at uncovering a potential performance bias when the train and test set have different class distribution, i.e. imbalanced and balanced respectively. Note that in this case, the Balanced Accuracy and Accuracy metrics are strictly

equivalent as we are evaluating performance on a balanced hold-out set. We therefore only report accuracy.

2.3.5. Significance testing on imbalanced data

We additionally evaluated the statistical significance of classification metrics across a range of imbalance ratios. Significance was computed from permutation tests with 100 permutations at $p < 0.01$, using an SVM trained on 1000 samples (at perfect balance) of synthetic data with different amounts of overlap between classes, namely: identical distributions (impossible classification problem), a distance of 1 between the means of the class distributions (difficult classification problem) and a distance of 3 between the two classes (easy classification problem) (Ojala and Garriga, 2010).

2.4. Brain data

To extend the analysis from a controlled environment with synthetic data towards a realistic setting with neuroimaging datasets, we ran experiments on publicly available EEG, MEG and fMRI datasets.

2.4.1. EEG – Motor Movement/Imagery Dataset

The publicly available EEG Motor Movement/Imagery Dataset (Goldberger et al., 2000; Schalk et al., 2004) consists of 64-channel EEG recordings of 109 subjects at 160Hz. While the dataset contains several tasks related to motor movement, only baseline resting-state runs were used, in this way creating a binary classification task between the eyes-open and eyes-closed conditions. Each recording has 1 minute of resting-state data which was segmented into 5-seconds epochs. As a result, 24 epochs per subject were extracted, half of them being eyes-closed and the others eyes-open. The effect of these conditions on neural oscillations is well studied, and consists of an increase in alpha power (8-12Hz) in posterior regions of the brain during eyes closed (Adrian and Matthews, 1934; Barry et al., 2007), compared to the eyes open condition. We computed alpha power (8-12Hz) from the power spectral density (PSD) obtained using the multi-taper method. To restrict the features to the visual cortex, only parieto-occipital electrodes (17 out of the 64) were used. The total sample size of this dataset was $109 \text{ subjects} \times 2 \text{ conditions} \times 12 \text{ trials} = 2616 \text{ samples}$ at perfect balance.

A second analysis was carried out on this dataset to study the relationship between electrode locations and performance scores along 3 different imbalance ratios (0.1, 0.5 and 0.9). As topographic differences are the main focus of this experiment, the sensors were not constrained to be from parieto-occipital regions. Similarly, placing the emphasis on performance as a function of spatial location, single-channel, single-feature SVM classifiers were trained (in contrast to multi-feature classification in the previous analysis). The motivation behind this analysis is to tentatively illustrate that with increasing class imbalance, classifiers lose focus on the areas whose data best discriminate the classes, and merely predict the majority class.

2.4.2. MEG – Cam-CAN Dataset

We used the passive auditory/visual perception task out of the open access MEG dataset collected at the Cambridge Centre for Ageing and Neuroscience (Cam-CAN) (Shafto et al., 2014). The preprocessing steps for this dataset can be found in Taylor et al. (2017). The task consists of 2-min recordings during which subjects were presented with either visual checkerboards or auditory tones (in random order) 60 times each, with a second between each stimulus. We further processed the data by down-sampling to 500Hz and epoching into 800-millisecond trials with 150 milliseconds of signal before stimulus onset and 650 milliseconds after. The epochs were baseline-corrected before computing alpha power (12-30Hz) using the multi-taper method on the 650 milliseconds after onset. We excluded the magnetometers and averaged powers for the two gradiometers for each location. For this study, we randomly selected 20 subjects out of the 643 that are available in the repository, resulting in a sample size of $20 \text{ subjects} \times 60 \text{ stimuli} \times 2 \text{ stimulus types} = 2400$

samples at perfect balance. Classification was performed on the data of a single channel (Fig. 5b), which was selected by training separate classifiers for all channels and selecting the location with best performance.

2.4.3. fMRI – Haxby dataset

Extending the analysis to functional MRI, we chose to use the publicly available Haxby dataset (Haxby et al., 2001) as provided through the Nilearn package (Abraham et al., 2014). Here we only explored recordings from subject 1 in a whole-brain voxel-wise classification paradigm. This allowed us to explore class imbalance in very high-dimensional feature spaces (39912 voxels). This dataset contains recordings of BOLD activity from individuals viewing images from different object categories. Stimuli were presented for 500ms with inter-stimulus intervals of 1500ms as part of a one-back repetition detection task. Here we trained four classifiers to predict the viewed object category from voxel-wise BOLD activity, limiting the object categories to faces and houses in order to have a binary classification task. This resulted in a sample size of $108 \text{ faces} + 108 \text{ houses} = 216$. Furthermore, this analysis serves as an example of within-subject classification as compared to the multi-subject setting in the previous tasks.

2.5. Data and code availability

The scripts, notebooks and pipeline used in this study are open-source under the MIT licence. The code is available on GitHub for further explorations. Our experiment pipeline is not limited to the datasets explored in this study and can easily be used to explore other datasets. The open-source repository can be found at: <https://github.com/thecocolab/data-imbalance>.

The code was developed using Python and its rich ecosystem for scientific computing. To process brain data we used MNE-Python (Gramfort et al., 2013), and machine learning algorithms and metrics came from scikit-learn (Pedregosa et al., 2011). Visualization was done with matplotlib (Hunter, 2007) and seaborn (Waskom, 2021).

The synthetic data used in this study can be generated using the open-source code we provide. The EEG Motor Movement/Imagery Dataset is publicly available and can be downloaded here: <https://physionet.org/content/eegmmidb/1.0.0/>. The Cam-CAN dataset can be accessed upon request at <https://camcan-archive.mrc-cbu.cam.ac.uk/dataaccess/> and the pipeline for preprocessing and loading the data is available at <https://github.com/arthurdehgan/camcan>.

3. Results

It is important to note here that it is not possible to compare absolute scores of metrics to one another. While for example AUC scores are generally higher than BAcc scores, this does not mean that the classifier performs better when using AUC. The performance metric does not influence the classifier's behavior. In fact, we used the same fitted instance of a classifier to evaluate all performance metrics. The only exception here is that Acc and BAcc are mathematically equivalent in the case of perfectly balanced data. In the following we focus mostly on the shape of the graph across imbalance ratios in order to compare performance metrics.

3.1. Simulated data

To demonstrate the effect of data imbalance on different performance metrics, algorithms and cross-validation techniques, we simulated three binary-class datasets, as described in subsection 2.1 and Fig. 2a, f, k.

3.1.1. Impossible classification

In the first scenario, the binary classification was performed on data from the same distribution (Fig. 2a-e), representing an impossible classification task. Due to the nature of the data, we expected performance

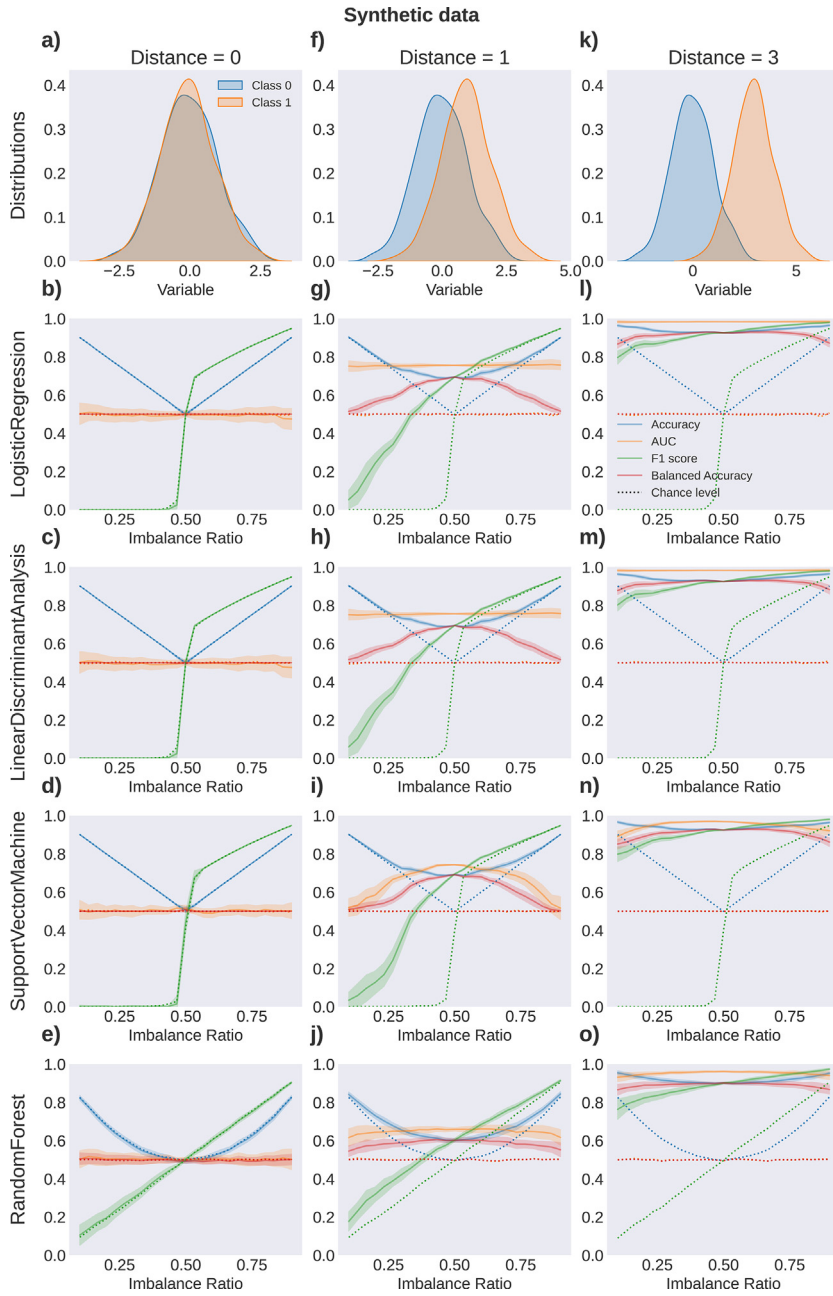


Fig. 2. Effect of data imbalance on different performance metrics and algorithms using simulated data. This figure summarizes results from three different synthetic datasets (see class distributions in the first row): column 1 (a–e): impossible classification task, all data comes from the same distribution, column 2 (f–j): strongly overlapping class distributions and column 3 (k–o): slightly overlapping class distributions. Rows 2 to 4 correspond to the different classification algorithms: Logistic Regression (b,g,l), Linear Discriminant Analysis (c,h,m), Support Vector Machine (d,i,n) and Random Forest (e,j,o). We evaluated Accuracy (blue), AUC (orange), F1 score (green), and BAcc (red). Solid lines show the performance over different class imbalance ratios, averaged over 10 initializations. Colored areas represent the respective standard deviation. Dashed lines indicate the average performance over 100 random permutations (i.e. chance level) for every performance metric. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

values to stay at the chance level, which we estimated using random permutations (Section 2.2). This hypothesis was confirmed by our results, with all performance metrics staying close to their respective chance level (not the probabilistic chance level of 0.5; it varies as a function of data imbalance and metric etc.) and showing only minimal variation across repetitions. We herein first describe the effect of different performance metrics and then describe the difference between classification algorithms. AUC and BAcc showed identical behavior, staying consistently at a chance level of 0.5 across all levels of class imbalance. In contrast, Accuracy scores and the respective chance level increased towards both extremes of data imbalance and reached minimal values at the point of perfect data balance. More specifically, the Accuracy score consistently reflected the proportion of the majority class in the imbalanced data (i.e. reaching a score of 90% for imbalance ratios of 9:1 and 1:9). The F1 score exhibited a steep increase from 0 towards 1 at the point where the data was balanced and continuously approached 1 with further increasing data imbalance.

Conspicuously, while Acc, AUC and BAcc show a symmetric pattern towards undersampling either class, the F1 score approaches 0 for increased undersampling of one class and 1 for the other. This behavior stems from F1 defining one class as positive and the other as negative, while the other evaluated metrics do not differentiate between the classes. F1 is a combined measure of the fraction of true positives among positive predictions (precision) and true positives among positive samples (recall). If the majority of samples are in the positive class, a classifier which always predicts the positive class will have precision and recall scores close to 1 and therefore high F1. On the other hand, an overrepresented negative class easily causes classifiers to always predict “negative”, resulting in an F1 score of 0. See Sibilini et al. (2020) for a version of F1 adapted to imbalanced data.

It is noteworthy that the behavior of all above-described metrics was most similar between LR, LDA, and SVM (linear and RBF kernel), but varied in RF. Compared to the other algorithms, accuracy using RF exhibited a slower increase towards the extreme imbalance and stayed

closer to 50% (i.e. the expected Accuracy score for classification between two datasets drawn from the same distribution) for a larger range of data imbalance ratios. In contrast to the steep increase of F1 described above, the F1 score in RF increased linearly over all levels of imbalance (Fig. 2b-e).

3.1.2. Difficult classification

In the second scenario, binary classification was performed on data from two overlapping Gaussian distributions with a distance of one (Fig. 2f), representing a difficult classification task. Due to the nature of this dataset we expected performance scores to reach above chance level. Over all levels of data imbalance, AUC remained consistent and above chance level in LR and LDA. Only in SVM and RF did the AUC decrease to chance level on both ends of increased data imbalance, with the strongest decrease in SVM. Accuracy exhibited a similar behavior over all four classification algorithms. Similarly to the first scenario, the Accuracy scores reached maximal scores on both ends of data imbalance. Despite the decrease of Acc towards the point of maximally balanced data, the distance between Acc and its chance level was increasing. The maximum distance between Acc and corresponding chance level was reached at the point of perfect class balance, indicating that the classifier was best at learning the structure in the dataset at this point. In contrast to Acc, BAcc reached its maximum at maximal data balance and dropped to chance level, i.e. 50%, at both ends of class imbalance. Again, the RF showed more stable BAcc scores over the range of class imbalance. The F1 score successively increased from 0 to 1 and showed similar but dampened behavior compared to the results from the identical class distributions (Fig. 2g-j). Results of the linear SVM model (Supp. Fig. 2) were similar to results of the SVM with RBF kernel. The only difference was in AUC, where we saw a flatter score around the balanced data regime.

3.1.3. Easy classification

In the third scenario, binary classification was performed on data from two Gaussian distributions with a distance of 3 (Fig. 2k), which is an example of an easy classification task. Due to the nature of the data, we expected performance values to reach high levels and to be less influenced by data imbalance. As expected, all performance metrics reached good classification scores which were less sensitive to data imbalance. While the general behavior was similar to the ones observed in the two previous experiments, the classifiers did a better job of learning structure from the data even in cases of imbalanced classes (Fig. 2l-o).

3.2. Secondary analysis

3.2.1. Hyperparameter tuning

We explored the effect on robustness when tuning hyperparameters related to class imbalance of LR, SVM and RF. However, we only found improved robustness towards imbalance for LR by weighting samples inversely proportional to class frequency. While enabling this re-weighting scheme led to a stable BAcc score across imbalance ratios for LR (Fig. 3a), SVM, and RF remained vulnerable to class imbalance (Fig. 3b, c). We examined a second hyperparameter for the Random Forest, namely the minimum weight fraction (MWFL) to generate a leaf node (Fig. 3c). Increasing this hyperparameter beyond its default value of zero led to a general improvement in BAcc (over the default hyperparameter set) for balanced data, and a decrease towards regions of extreme class imbalance. Therefore, we are not able to report improvements in robustness from this hyperparameter for the Random Forest.

3.2.2. Sample size

We did not find any notable effect of sample size ($N=300$, $N=1000$, $N=3000$) on overall SVM classification accuracy nor robustness to class

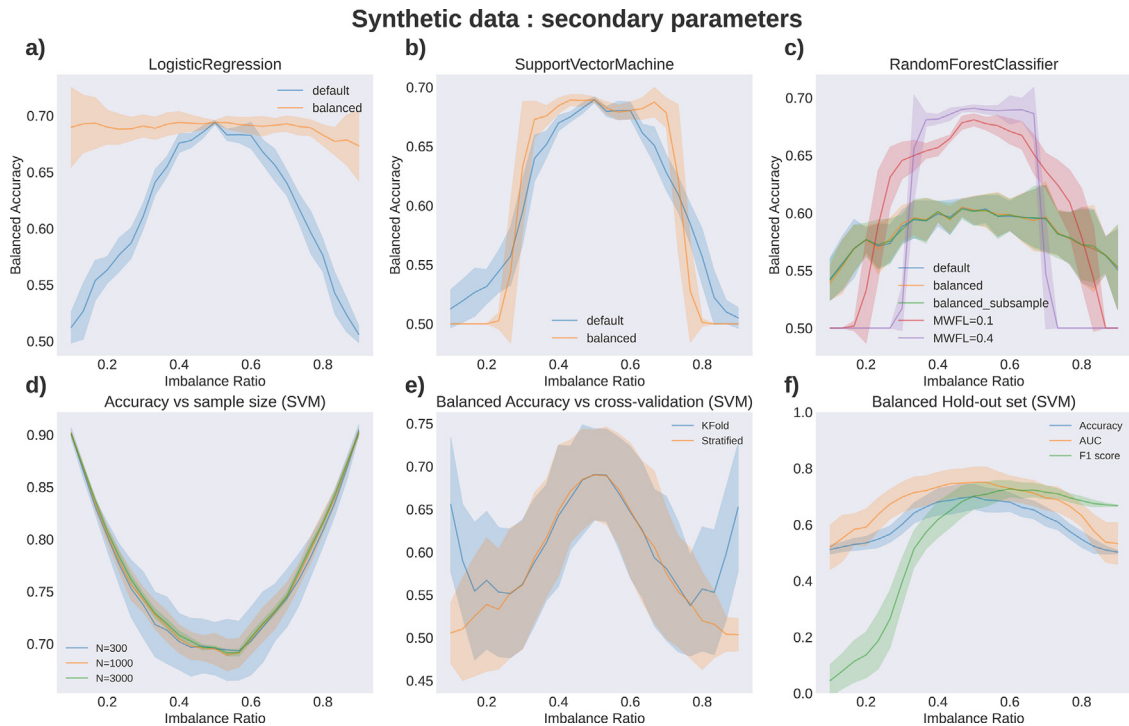


Fig. 3. These examples follow the same principle as previous analysis (Fig. 2), expanding on the effect of secondary parameters and hyperparameter tuning on the robustness to imbalance. (a-c) Exploration of the effect of hyperparameters on classifier robustness against imbalance. Each line represents a certain hyperparameter setting (Section 2.3.1) d) The effect of sample size on robustness to class imbalance using an SVM and Acc. e) The impact of cross-validation scheme on robustness to class imbalance (BAcc and SVM). f) Performance metrics on a balanced hold-out set using an SVM trained on different ratios of imbalance. BAcc is not shown here as it is equivalent to Acc for balanced data. Solid lines show the performance over different class imbalance ratios, averaged over 10 initializations. Colored areas represent the respective standard deviation.

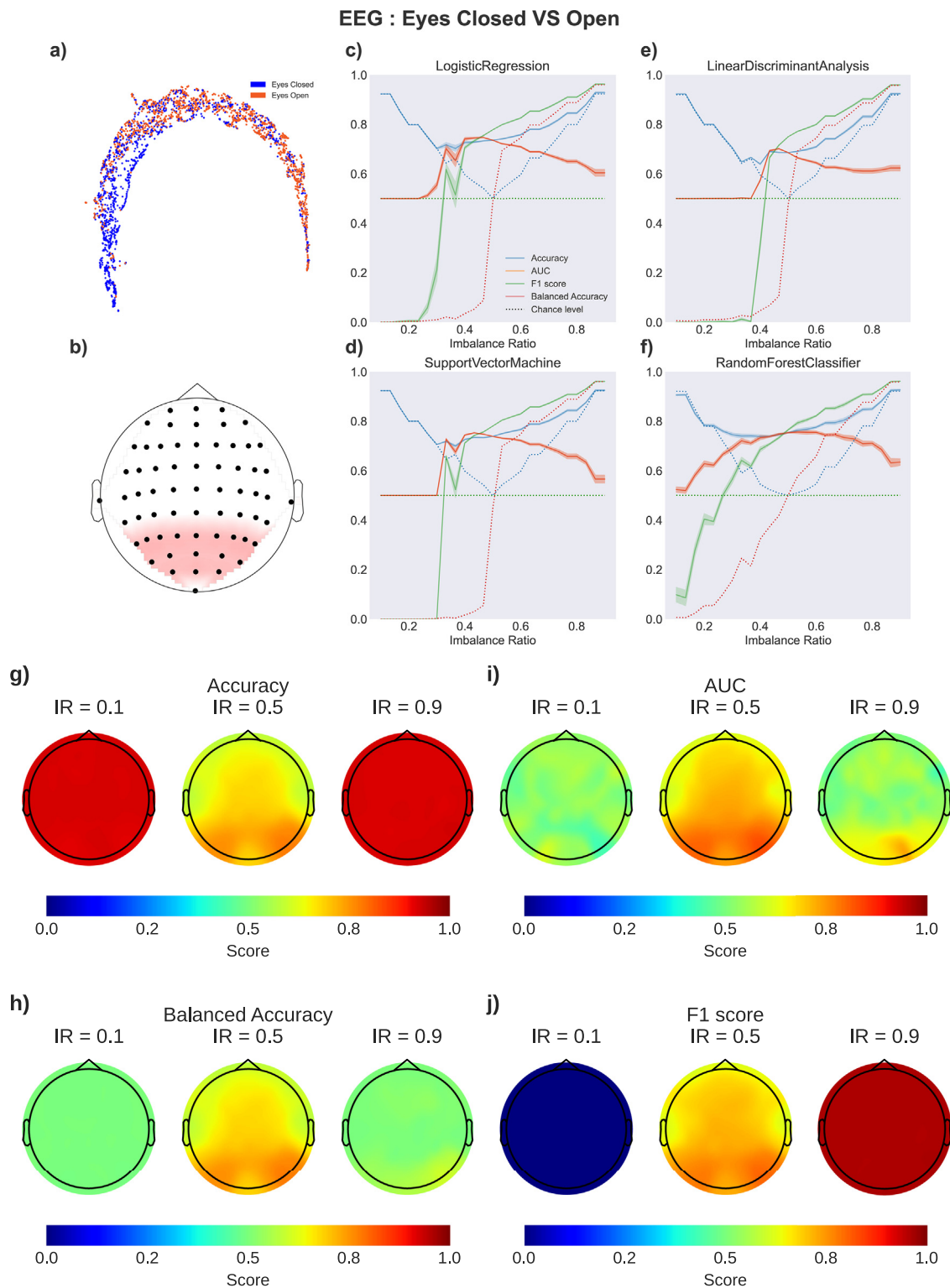


Fig. 4. Effect of data imbalance on different performance metrics and algorithms for EEG data classification. **a)** 2D projection of the 17-dimensional input space using the UMAP algorithm (McInnes et al., 2018). Each dot represents one sample mapped onto two UMAP components (x- and y-axis). This illustrates the amount of overlap (related to classification difficulty) between the eyes closed (blue) and eyes open (orange) classes. **b)** The parieto-occipital region of interest (ROI) was used as features. The effect of data imbalance using **c)** Logistic Regression; **d)** Support Vector Machine; **e)** Linear Discriminant Analysis, and **f)** Random Forest. We evaluated Accuracy (blue), AUC (orange), F1 score (green), and BAcc (red). Solid lines indicate performance across imbalance ratios averaged over 10 random initializations. Colored areas represent the respective standard deviation. Dotted lines indicate the average score across 100 random permutations of class labels (i.e. data-driven chance level). Subfigures **g-j)** highlight the effect between data imbalance and sensor location across performance metrics for single-channel, single-feature classification between eyes-open and eyes-closed EEG using an SVM. IR indicates the imbalance ratio used for each topographical map. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

imbalance. However, as one would expect, the variance across random initializations of the cross-validation splits decreased with the number of training samples (Fig. 3d). This result validates our procedure of undersampling to achieve class imbalance, which leads to a decrease in sample size towards higher degrees of imbalance. As here we show that a change in sample size only affects the variance of the results, the shape of the performance curves can be compared across imbalance ratios.

3.2.3. Cross-Validation

Comparing different cross-validation algorithms revealed that K-Fold is more sensitive to imbalanced data than Stratified K-Fold (Fig. 3e). While both procedures led to similar Balanced Accuracy scores with balanced classes, K-Fold cross-validation without stratification showed an increase in Balanced Accuracy towards extremely imbalanced datasets. This likely stems from validation splits containing only a single class, which Balanced Accuracy does not account for. While a classifier voting for the majority class has a Balanced Accuracy of 0.5 on data containing two classes, its score will be 1 on validation splits that, by chance, only contain one class. The likelihood of this is higher for low sample sizes and is completely resolved by using Stratified K-Fold cross-validation.

3.2.4. Balanced hold-out set

So far, all experiments were evaluated using cross-validation with imbalanced validation splits, replicating the class distribution of the training set. Fig. 3f shows SVM classification scores on a balanced hold-out set after being trained on increasingly more imbalanced training data. While AUC and F1 scores are largely in line with previous results (Fig. 2i) using imbalanced cross-validation splits, Acc now reflects the

previous behavior of BAcc, i.e. dropping towards the random baseline of 50% towards the extremes of class imbalance. Note that we do not report BAcc here as it is equivalent to Acc on balanced data.

3.2.5. Significance testing on imbalanced data

Supplementary Fig. 1 depicts significance scores across imbalance ratios and levels of difficulty of the classification problem. Generally we found that for the impossible classification task (i.e. identical class distributions; Supp. Fig. 1a), none of the scores were significant at $p < 0.01$, as in this case, permuting labels does not remove any structure from the data. For the difficult classification task b), however, we found a range of statistically significant classification scores around perfect class balance. Scores were not significant towards the extremes of class imbalance. This behavior was shared among all the classification metrics we evaluated. In the third task—easy classification (Supp. Fig. 1c)—all classification scores were found to be significantly above chance level for all metrics, which highlights the classifier's ability to learn structure in the data even for extreme class imbalance, which is even more pronounced for easier classification tasks. Note that the results in Supp. Fig. 1 differ from the results presented in Fig. 2d,i,n because — even though the experimental setup was the same — when testing for statistical significance we are limited to a single repetition of the analysis, while results in Fig. 2 are averaged over 10 random seeds.

3.3. Results on EEG data

We performed two experiments using the EEG dataset with classification between eyes-open versus eyes-closed during resting

MEG Cam-CAN : Visual VS Auditory stimulation

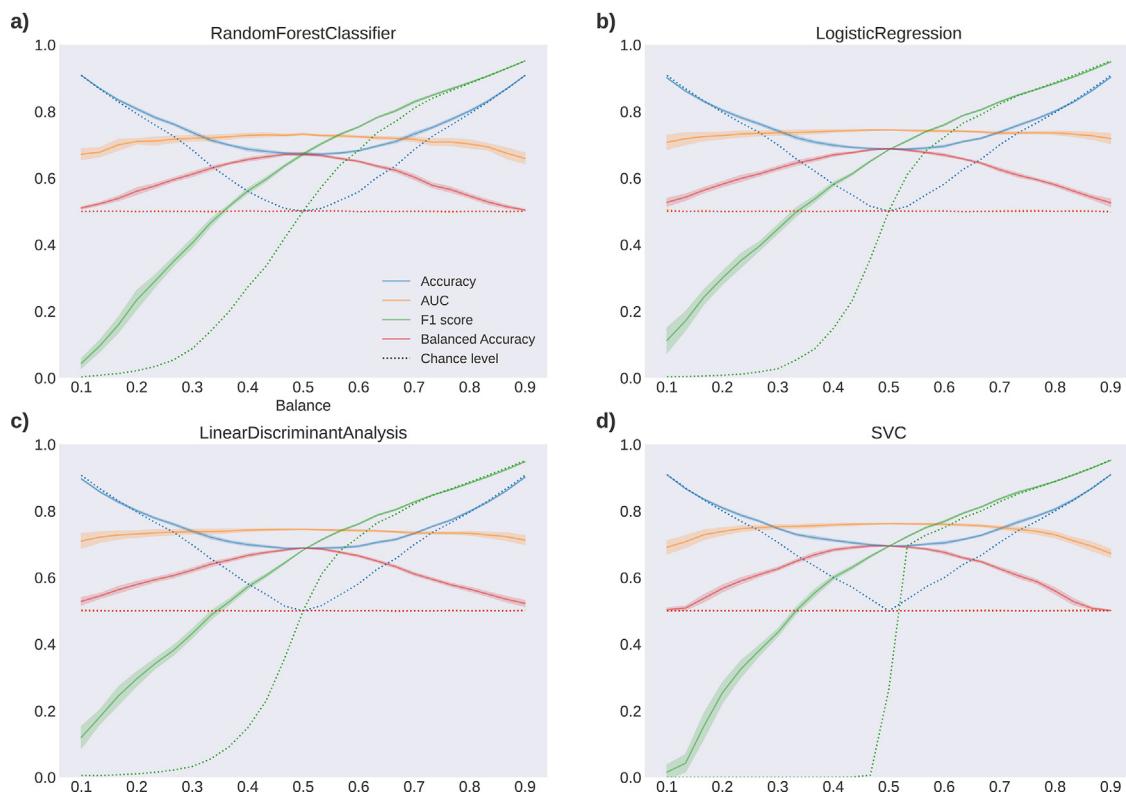


Fig. 5. Effect of data imbalance on different performance metrics and algorithms on classification of MEG data. The effect of class imbalance using a) Logistic Regression; b) Random Forest; c) Linear Discriminant Analysis; d) Support Vector Machine. We evaluated Accuracy (blue), AUC (orange), F1 score (green), and BAcc (red). Solid lines indicate the performance over different class imbalances, averaged over 10 initializations. Colored areas represent the respective standard deviation. Dotted lines indicate the performance of 100 random permutations (i.e. chance level) for every performance metric. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

(Section 2.4.1). We herein first describe the results of the multi-feature classification using parieto-occipital electrodes and then present the results of the channel-wise classification.

As previously shown in the simulated data, Acc increased with increasing data imbalance, reflecting the proportion of the majority class. In contrast to this, BAcc approached chance level with increased imbalance and reached its maximum with maximally balanced data. In line with the simulated data, the F1 score increased abruptly at optimal class balance. AUC was stable over a wide range of class imbalance. While SVM, LR, and LDA showed similar behavior (i.e. being equally sensitive to data imbalance), the performance metrics using RF exhibited more stability over different levels of imbalance (Fig. 4c-f).

In contrast to multi-channel classification, single-channel decoding performance is commonly used to localize changes between two conditions and allows to attribute larger changes to channels with higher decoding performance (Fig. 4a, b). Here we show how class imbalance may lead to misinterpretation of such results and the loss of structure related to decoding performance. At optimal class balance, highest decoding performance was found in parieto-occipital regions of the brain, which is in line with the literature and allows for interpretation of the results (Fig. 4g-j, $IR = 0.5$). Towards the extremes of data imbalance, this structure is lost and we find uniform decoding performance across the brain. This effect is most prominent using Acc, BAcc, and F1, while AUC retains some variations across channels (Fig. 4g-j).

3.4. Results on MEG data

As described in Section 2.4.2, classification was performed between auditory and visual stimuli, using all sensor locations as features and averaged power across the two MEG gradiometers for each location. Fig. 5 shows the results of the multi-feature classification using different models. All models seem to perform similarly and have similar reaction to imbalance with all of the studied performance metrics. The only notable difference in baseline chance level computations appears for the SVC classifier where F1 score raises sharply when reaching the balanced data ratio (Fig. 5). In line with the simulated and EEG datasets, BAcc approached chance level with increased imbalance and reached its maximum with perfectly balanced data. While AUC was stable across a wide range of imbalance ratios for LR, LDA, and RF, we observed comparable behavior to BAcc for SVM (i.e. approaching chance level towards more data imbalance). While SVM, LR, and LDA showed similar behavior (i.e. being equally sensitive to data imbalance), the performance metrics in RF exhibited slightly more stability over different levels of imbalance. This is in line with the aforementioned results using the simulated and EEG datasets.

3.5. Results on fMRI data

Fig. 6 shows the results of our analysis in a high-dimensional classification setting using whole-brain voxel-wise classification. In line with previous observations using the same dataset (cf. Nilearn tutorials (Abraham et al., 2014)), all tested classifiers performed extremely well in the balanced data condition often exceeding 95% decoding accuracy. Importantly, the results we observed here with regards to the behavior of the four performance metrics across the classifiers are in line with the previous results. As expected from the classification task (face vs. house stimulus), the Logistic Regression weights peak in the fusiform face area (FFA), confirming that the high decoding accuracy stems from meaningful learned representations (Fig. 6a).

4. Discussion

The present work shows that the implementation of classification on imbalanced data is feasible, though it demands certain important considerations. Our approach demonstrates how one needs to be mindful of class imbalance when choosing a classifier, an evaluation metric and a

cross-validation scheme (Krawczyk, 2016). Here, we sought to provide a didactic technical note on this question using a combination of simulated data, electrophysiological brain signals and fMRI recordings. Concretely, we quantified the behavior of commonly used classifiers, performance metrics, and cross-validation approaches across varying levels of data imbalance. An exhaustive exploration of all available techniques that have been proposed to tackle data imbalance is beyond the scope of this study. Instead, we chose to focus on machine learning tools and metrics that are often used within the neuroscience community. In line with this, the methods we address—and the open-source pipelines and notebooks we provide—all use the scikit-learn library.

Taken together, our observations support the idea that classification on moderately imbalanced data is feasible, as long as appropriate classifiers and performance metrics are employed. More specifically, by systematically manipulating the degree of data imbalance, we illustrated and quantified several key effects. First, we confirmed the tendency of classifiers to resort to blindly voting for the majority class as data imbalance was accentuated. When assessed with the widely used Accuracy measure, this behavior was associated with an artificial improvement in the model's classification performance. AUC and BAcc were more robust to the increase in imbalance, and are therefore more appropriate under these circumstances. Moreover, our data confirms that Random Forest is more robust when handling imbalanced data, compared to other commonly used algorithms, such as LR, LDA, and SVM—especially when using class-weighting hyperparameter optimization. This result is expected, but our analyses quantify this for a wide range of imbalances and illustrates the effect with simulated data as well as EEG and MEG recordings. We also found that the balancing hyperparameter can be used to improve LR's robustness to data imbalance.

Our study also highlights an important caveat concerning the use of permutation tests on imbalanced data. Permutation tests allow data-driven computation of the chance level, and from this chance level, they provide an estimate of statistical significance. However, because chance levels can be much greater than 50% in imbalanced data (for binary classification problems; Fig. 2), a simple reporting in these cases of the Accuracy and of its statistical significance can artificially inflate the importance of the classification result. For example, for a 0.2 imbalance ratio, a statistically significant Accuracy of 82% can appear to be an outstanding result, when in reality, the chance level is 80%, so this is arguably comparable in importance to a statistically significant Accuracy of 52%. This scenario underlines the importance of reporting the chance level alongside performance metrics when carrying out permutation testing on imbalanced data (Combrisson and Jerbi, 2015).

The present study complements a wide array of insightful investigations that have explored the pitfalls and potential solutions for supervised learning with imbalanced data (Dubey et al., 2014; Graa and Rekik, 2019; Japkowicz and Stephen, 2002; Jeni et al., 2013; Kamalov and Denisov, 2020; Krawczyk, 2016; Prati et al., 2015; Straube and Krell, 2014; Sun et al., 2009; Tan et al., 2007; Thabtah et al., 2020; Varoquaux and Colliot, 2022; Wang et al., 2017; Wardhani et al., 2019). By contrast to some of the previous studies, the present work focuses on insights that are directly relevant to researchers in neuroimaging using standard tools, such as those available through the scikit-learn library.

We would like to emphasize that while the examples provided here represent a specific instance of classification problems that can be encountered in the field of neuroscience, our analysis serves as an illustrative placeholder for many types of problems including but not limited to within-subject or between-subject classification, low- or high-dimensional brain imaging data, or analysis of behavioral data. Importantly, we encourage the reader to follow the same procedure to explore the imbalance question in their own data. As a matter of fact, the code and Jupyter notebooks (<https://github.com/thecocolab/data-imbalance>) were designed so the figures could be easily replicable, and allow users to extend the investigations to a wider range of metrics and methods, in a collaborative and open science perspective.

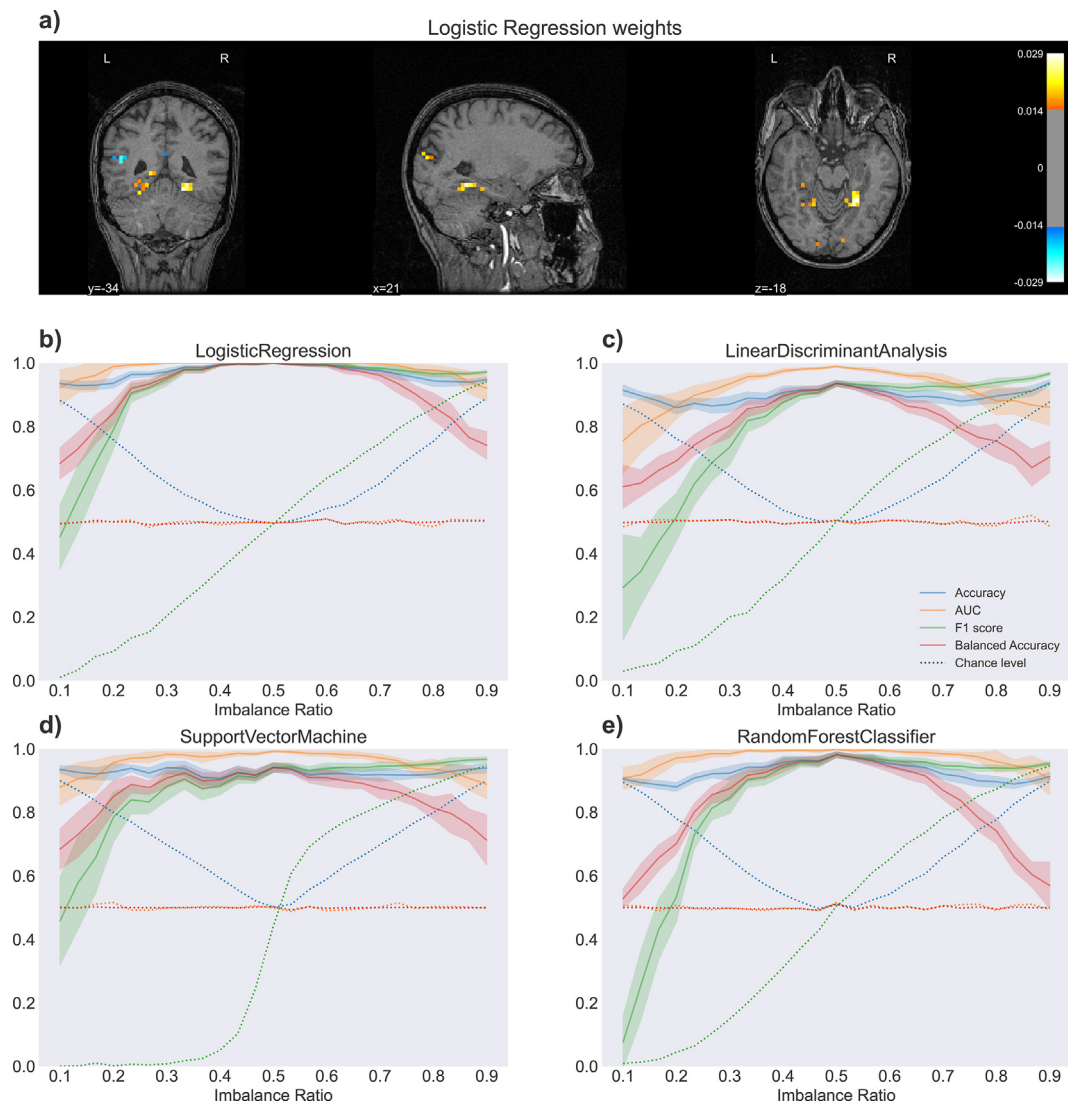


Fig. 6. Effect of data imbalance on different performance metrics, algorithms, and cross-validation techniques on classification of fMRI data. **a)** Illustration of Logistic Regression weights trained on classifying face and house stimuli. The effect of class imbalance using **b)** Logistic Regression; **c)** Linear Discriminant Analysis; **d)** Support Vector Machine; **e)** Random Forest. We evaluated Accuracy (blue), AUC (orange), F1 score (green), and BAcc (red). Solid lines indicate the performance over different class imbalances, averaged over 10 initializations. Colored areas represent the respective standard deviation. Dotted lines indicate the performance of 100 random permutations (i.e. chance level) for every performance metric. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Recommendations In discussing recommendations and best practices for handling data imbalance, it is important to note that the utility or suitability of a metric should always be determined in relation to the specific problem being tackled. The various types of machine learning problems differ among other things in the type of error one seeks to minimize—this will in turn determine the appropriate classification paradigm in which to operate. The suitable paradigm depends on whether one’s aim is to minimize overall classification error (Lee et al., 2021), type I error (false positive) (Van De Ruit and Grey, 2019) or type II error (false negative) (Abdelhamid et al., 2020). Of these three, the design of the present study and the findings we report relate primarily to paradigms seeking to minimize overall classification error. Our analysis confirms that BAcc (and the related AUC metric) are more robust to data imbalance than the common Accuracy metric. Note that the Precision and Recall metrics are recommended when dealing with type I and type II errors respectively. Sampling techniques (i.e. over- and undersampling) are known to be helpful in most paradigms and evidence suggests that they work well in combination with certain classifiers (Feng et al., 2020).

Based on the observations reported in this study, evaluating the performance metrics on a balanced holdout set (Fig. 3f) allows for an unbiased evaluation of classification performance even when using metrics vulnerable to class imbalance, such as Acc. We therefore recommend using several evaluation metrics, e.g. BAcc, with Stratified K-Fold cross-validation and standard Accuracy on a balanced holdout set. As BAcc is equivalent to Acc for balanced data, it retains Acc’s greatest advantage, namely its intuitiveness and ease of interpretation. We additionally argue that BAcc results in more intuitive performance evaluation for imbalanced data, as it combines performances of individual classes with equal weight. Accuracy on the other hand combines class performances with a strong bias towards the majority class.

Deep learning, an advanced type of machine learning used for a large variety of classification tasks, is not immune to data imbalance (Buda et al., 2018). Deep learning models learn by backpropagating gradients through the model. In class-imbalanced scenarios, the majority class dominates the net gradient that is responsible for updating the model’s weights, which reduces the error of the majority group quickly during early iterations. However, oftentimes it simultaneously increases

the error of the minority group. As a result, the neural network struggles to learn the decision boundary for the problem (Anand et al., 1993). Common approaches used to overcome data imbalance when training deep neural networks include under-/oversampling, data augmentation, the use of a class-weighted loss function (i.e. higher penalty for errors made on the minority class) and output thresholding (Johnson and Khoshgoftaar, 2019).

Overall, given our results, we make the following seven recommendations for machine learning in neuroscience:

1. **Know your problem.** The right performance metric to use is determined first and foremost by the specific research question at hand. Throughout this study, we have focused largely on the type of question that requires a minimization of the overall classification error. However, for some problems, it may be more important to prioritize the minimization of either type I or type II error. In these cases, other performance metrics may be more relevant than BAcc or AUC, namely Precision and Recall. For example, a classifier used to discriminate biological sex based on brain activity would likely seek to minimize overall classification error, and BAcc or AUC would then be recommended to correct for any class imbalance in the data. In contrast, a classifier used to detect malignant tumours in brain imaging would instead benefit most from Precision and Recall metrics, since false negatives (type II error) would need to be minimized above all else. Therefore, consideration of the nature of the problem at hand—specifically, of the type of error to be minimized—is essential to selecting the most appropriate performance metric. With this in mind, the remaining recommendations offered here apply specifically to problems requiring the minimization of classification error.
2. **Use Balanced Accuracy (BAcc).** BAcc has been largely underexploited in neuroscience research. Given (i) its superior robustness to imbalance, (ii) the fact that it simply reduces to Accuracy for balanced datasets and (iii) its applicability to both binary and multi-class (Grandini et al., 2020) datasets, we recommend the routine use of BAcc, rather than the commonly used Acc, as a default for neuroscience machine learning applications where overall classification error should be minimized. To maximize the interpretability of classification results however, it is worth looking at multiple performance metrics (e.g. BAcc and AUC). If the classifier is purely used to decide if a feature captures a difference between two conditions (as commonly done in the field of neuroscience), AUC combined with significance testing serves as a powerful tool.
3. **Use ensemble methods.** If data imbalance cannot be avoided, we recommend the use of classifier families that provide additional robustness. In line with previous recommendations, ensemble methods such as Random Forests are less sensitive to data imbalance and provide a set of hyperparameters that can be optimized for the classification of imbalanced data. Further, ensemble methods are generally known to improve robustness (Dietterich, 2000; Sagi and Rokach, 2018). However, the complexity of some of these algorithms might not fit well onto very simple classification tasks, which may affect the robustness.
4. **Use a balanced hold-out test set.** When working with imbalanced data, the true overall classification error of the trained classifier can be assessed simply by testing it on a balanced hold-out test set. However, we further suggest also evaluating classifier performance on a test set that reflects the class distribution of the specific problem in the wild (or of the training set, if the class distribution in the real-life setting cannot be estimated). The difference in score between these two test sets can additionally help interpret the performance of the classifier. Beyond that, analysing confusion matrices and their derived metrics helps with understanding the behavior of classifiers in more detail and shed light on class imbalance related biases. For an illustration of this see Hahn et al. (2013).
5. **Use Stratified K-fold for cross-validation.** K-fold cross validation is highly sensitive to data imbalance such that in extreme cases, it can even become impossible to perform classification (i.e. if one fold contains only a single class). We therefore strongly recommend the use of Stratified K-fold, which maintains the imbalance ratio within each of the selected folds.
6. **Report statistical significance AND chance level.** Without the corresponding chance level, performance metrics can easily be misinterpreted. Especially for performance metrics like Acc, chance level fluctuates widely with data imbalance. Thus, performance scores should always be reported accompanied by—and should be interpreted in the light of—the associated chance level. The chance level can be estimated in a data-driven approach using permutation tests. This recommendation also applies if statistical tests were performed and performance scores reached significance over random permutations.
7. **Use hyperparameters.** For many of its classifiers, scikit-learn provides hyperparameter options specifically designed for dealing with imbalanced data; these should be routinely exploited. Our study highlights the potential utility of this step and encourages the reader to consider hyperparameters for an optimal performance. An extensive tutorial on hyperparameter selection is beyond the scope of this work. There is substantial useful literature on hyperparameter selection for further reading (Andonie, 2019; Glaser et al., 2020; Hosseini et al., 2020; Lemm et al., 2011; Luo, 2016; skl, 0000; Yang and Shami, 2020).

Limitations and perspectives This study's results need to be interpreted in the light of several limitations. First, many approaches for handling imbalanced data have been proposed (Haixiang et al., 2017). In this study, we focused on families of models and performance metrics that are easily accessible and widely used in the neuroscience community, in particular through the scikit-learn library. Specifically, we focused on four popular classifiers and four standard evaluation metrics, as well as two cross-validation schemes. Reviewing or comparing all existing tools is beyond the intended goal of this paper, but the Python code we provide is open-source, which allows users to extend these investigations.

Second, the data we used consisted of simulated Gaussian data distributions, as well as open-access electrophysiological (MEG and EEG) and fMRI brain signals. We did not examine the impact of noise, though it could have been interesting to consider it, as it has been shown to interact with the performance of the classifier (Somasundaram and Reddy, 2017).

Third, we only briefly mention the option of balancing the data through over and under-sampling methods, and focused our investigation on evaluating the impact of data imbalance.

Fourth, the aim of this study was to explore the individual classifiers' sensitivities and relative changes of performance metrics over systematically imbalanced data. Therefore, the presented results do not allow a comparison of individual classifiers' absolute performance.

Fifth, our main focus was to explore the effect of data imbalance on classification tasks. To this end, we explored the simplest form, namely a binary classification, where the manipulation of the imbalance ratio is straightforward. While creating imbalanced data for the multiclass case is more complicated, most of our findings extend to this type of classification—with some caveats. For instance, AUC is typically available only for binary classification, although extensions for multiclass classification exist (Hand and Till, 2001). In addition, some algorithms such as SVM and LR reframe multi-class classification as a separate one-vs-the-rest binary classification problem for each class; this effectively results in class imbalance. While other options exist, such as multinomial logistic loss (i.e. cross-entropy loss) for LR or one-vs-one classification for SVMs, we advise caution when using these algorithms in a multi-class setting. Notwithstanding, the main take-home message of this study—that BAcc is generally more robust than Acc—still applies for multiclass classification.

Conclusion In this study, we have illustrated the effect of imbalanced data on some of the most prominent classification algorithms and

performance metrics used in neuroscience. Among other things, one key take-home message is our suggestion to systematically use Balanced Accuracy over the widely used Accuracy metric, whenever the aim is to minimize overall classification error. In addition to its robustness to class imbalance, Balanced Accuracy collapses to standard Accuracy for balanced datasets and is readily extendable to multiclass problems. More generally, we hope that the recommendations and red flags reported here—using simulations and real brain data—will strengthen good practices for the application of supervised ML in neuroscience, and increase awareness, especially among new-comers to the field. Lastly, by providing open-source code and well-documented pipelines, we hope that others will further explore this question with a wider variety of parameters, classifiers, and different types of classification problems.

Declaration of Competing Interest

The authors declare no conflict of interest.

Credit authorship contribution statement

Philipp Thölke: Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Yorguin-Jose Mantilla-Ramos:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Hamza Abdelhedi:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Charlotte Maschke:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Arthur Dehgan:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Yann Harel:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Anirudha Kemptur:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Loubna Mekki Berrada:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Myriam Sahraoui:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Tammy Young:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Antoine Bellemare Pépin:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Clara El Khantour:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Mathieu Landry:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Annalisa Pascarella:** Methodology, Writing – review & editing. **Vanessa Hadid:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Etienne Combrisson:** Methodology, Writing – review & editing. **Jordan O’Byrne:** Conceptualization, Methodology, Formal analysis, Writing – review & editing. **Karim Jerbi:** Conceptualization, Methodology, Formal analysis, Writing – review & editing.

Data availability

All code and data are openly accessible. See the Data and Code Availability section.

Acknowledgements

KJ was supported by funding from the Canada Research Chairs program and a Discovery Grant from the Natural Sciences and Engineering Research Council of Canada (NSERC), and IVADO- Apogée fundamental research project grant. YH was supported by the Courtois-Neuromod Project. CM was supported by the FRQNT Strategic Clusters Program (2020-RS4-265502 - Centre UNIQUE - Union Neurosciences and Artificial Intelligence Quebec and the Canada First Research Excellence Fund and Fonds de recherche du Québec, awarded to the Healthy Brains, Healthy Lives initiative at McGill University. HA was supported by Mitacs Globalink Research Award. JOB was supported by the Canadian Institutes of Health Research (CIHR). AD was supported by Mila. AK was supported by Mitacs Graduate Fellowship and Courtois foundation. LMB was supported by IVADO. PT acknowledges support through a

scholarship from the Cognitive and Computational Neuroscience Laboratory (CoCo Lab) and Mila (Quebec Machine Learning Institute). TY was supported by Mitacs. YM was supported by Mitacs through the Mitacs Globalink Research Internship program. The authors wish to thank Kris, Mark, T. Hortons and Mr. and Mrs. Puffs for their support during the Hackathon/Paper Sprint events.

Supplementary material

Supplementary material associated with this article can be found, in the online version, at [10.1016/j.neuroimage.2023.120253](https://doi.org/10.1016/j.neuroimage.2023.120253)

References

- Abdelhamid, N., Padmavathy, A., Peebles, D., Thabtah, F., Goulder-Horobin, D., 2020. Data imbalance in autism pre-diagnosis classification systems: an experimental study. *J. Inf. Knowl. Manag.* 19 (01), 2040014.
- Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., Varoquaux, G., 2014. Machine learning for neuroimaging with scikit-learn. <https://www.frontiersin.org/articles/10.3389/fninf.2014.00014/full>.
- Adrian, E.D., Matthews, B.H., 1934. The berger rhythm: potential changes from the occipital lobes in man. *Brain* 57 (4), 355–385.
- Anand, R., Mehrotra, K., Mohan, C., Ranka, S., 1993. An improved algorithm for neural network classification of imbalanced training sets. *IEEE Trans. Neural Netw.* 4 (6), 962–969. doi:10.1109/72.286891.
- Andonie, R., 2019. Hyperparameter optimization in learning systems. *J. Membrane Comput.* 1 (4), 279–291. doi:10.1007/s41965-019-00023-0.
- Barry, R.J., Clarke, A.R., Johnstone, S.J., Magee, C.A., Rushby, J.A., 2007. Eeg differences between eyes-closed and eyes-open resting conditions. *Clin. Neurophysiol.* 118 (12), 2765–2773.
- Bode, S., Feuerriegel, D., Bennett, D., Alday, P.M., 2019. The decision decoding toolbox (ddtbox)—a multivariate pattern analysis toolbox for event-related potentials. *Neuroinformatics* 17 (1), 27–42.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45 (1), 5–32.
- Brodersen, K.H., Ong, C.S., Stephan, K.E., Buhmann, J.M., 2010. The balanced accuracy and its posterior distribution. In: 2010 20th International Conference on Pattern Recognition, pp. 3121–3124. doi:10.1109/ICPR.2010.764.
- Brodersen, K.H., Ong, C.S., Stephan, K.E., Buhmann, J.M., 2010. The balanced accuracy and its posterior distribution. In: 2010 20th international conference on pattern recognition. IEEE, pp. 3121–3124.
- Buchlak, Q.D., Esmaili, N., Leveque, J.-C., Bennett, C., Farrokhi, F., Piccardi, M., 2021. Machine learning applications to neuroimaging for glioma detection and classification: an artificial intelligence augmented systematic review. *J. Clin. Neurosci.* 89, 177–198.
- Buda, M., Maki, A., Mazurowski, M.A., 2018. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Netw.* 106, 249–259. doi:10.1016/j.neunet.2018.07.011.
- Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. Smote: synthetic minority over-sampling technique. *J. Artif. Intell. Res.* 16, 321–357.
- Chawla, N.V., Japkowicz, N., Kotcz, A., 2004. Special issue on learning from imbalanced data sets. *ACM SIGKDD Explorat. Newsletter* 6 (1), 1–6.
- Combrisson, E., Jerbi, K., 2015. Exceeding chance level by chance: the caveat of theoretical chance levels in brain signal classification and statistical assessment of decoding accuracy. *J. Neurosci. Methods* 250, 126–136. doi:10.1016/j.jneumeth.2015.01.010.
- Cutting-edge EEG Methods
- Cortes, C., Vapnik, V., 1995. Support-vector networks. *Mach. Learn.* 20 (3), 273–297.
- Cox, D.R., 1958. The regression analysis of binary sequences. *J. R. Stat. Soc.: Ser. B (Methodological)* 20 (2), 215–232.
- Das, A., Geisler, W.S., 2021. A method to integrate and classify normal distributions. *J. Vis.* 21 (10), 1–1.
- Dietterich, T.G., 2000. Ensemble methods in machine learning. In: *Multiple Classifier Systems*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 1–15.
- Dubey, R., Zhou, J., Wang, Y., Thompson, P.M., Ye, J., Initiative, A.D.N., et al., 2014. Analysis of sampling techniques for imbalanced data: an n= 648 adni study. *Neuroimage* 87, 220–241.
- Fahrenfort, J.J., Van Driel, J., Van Gaal, S., Olivers, C.N., 2018. From erps to mvpa using the amsterdam decoding and modeling toolbox (adam). *Front. Neurosci.* 12, 368.
- Feng, Y., Zhou, M., Tong, X., 2020. Imbalanced classification: an objective-oriented review. *arXiv preprint arXiv:2002.04592*.
- Fernández, A., García, S., Herrera, F., Chawla, N.V., 2018. Smote for learning from imbalanced data: progress and challenges, marking the 15-year anniversary. *J. Artif. Intell. Res.* 61, 863–905.
- Fong, R.C., Scheirer, W.J., Cox, D.D., 2018. Using human brain activity to guide machine learning. *Sci. Rep.* 8 (1), 1–10.
- Fukunaga, K., 1993. *Statistical Pattern Recognition*. In: *Handbook of pattern recognition and computer vision*. World Scientific, pp. 33–60.
- Gershman, S.J., Horvitz, E.J., Tenenbaum, J.B., 2015. Computational rationality: a converging paradigm for intelligence in brains, minds, and machines. *Science* 349 (6245), 273–278.
- Glaser, J.I., Benjamin, A.S., Chowdhury, R.H., Perich, M.G., Miller, L.E., Kording, K.P., 2020. Machine learning for neural decoding. *eNeuro* 7 (4).
- Glaser, J.I., Benjamin, A.S., Farhoodi, R., Kording, K.P., 2019. The roles of supervised machine learning in systems neuroscience. *Prog. Neurobiol.* 175, 126–137.

- Goldberger, A.L., Amaral, L.A., Glass, L., Hausdorff, J.M., Ivanov, P.C., Mark, R.G., Mietus, J.E., Moody, G.B., Peng, C.-K., Stanley, H.E., 2000. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *Circulation* 101 (23), e215–e220.
- Gong, M., 2021. A novel performance measure for machine learning classification. *Int. J. Manag. Inf. Technol. (IJMIT)* Vol 13.
- Graa, O., Reikik, I., 2019. Multi-view learning-based data proliferator for boosting classification using highly imbalanced classes. *J. Neurosci. Methods* 327, 108344.
- Gramfort, A., Luessi, M., Larson, E., Engemann, D.A., Strohmeier, D., Brodbeck, C., Goh, R., Jas, M., Brooks, T., Parkkonen, L., Hämäläinen, M.S., 2013. MEG And EEG data analysis with MNE-python. *Front. Neurosci.* 7 (267), 1–13. doi:10.3389/fnins.2013.00267.
- Grandini, M., Bagli, E., Visani, G., 2020. Metrics for multi-class classification: an overview. doi:10.48550/ARXIV.2008.05756.
- Grootswagers, T., Wardle, S.G., Glass, T.A., 2017. Decoding dynamic brain patterns from evoked responses: a tutorial on multivariate pattern analysis applied to time series neuroimaging data. *J. Cogn. Neurosci.* 29 (4), 677–697.
- Hahn, T., Marquand, A.F., Plichta, M.M., Ehli, A.-C., Schecklmann, M.W., Dresler, T., Jarczok, T.A., Eirich, E., Leonhard, C., Reif, A., Lesch, K.-P., Brammer, M.J., Mourao-Miranda, J., Fallgatter, A.J., 2013. A novel approach to probabilistic biomarker-based classification using functional near-infrared spectroscopy. *Hum. Brain Mapp.* 34 (5), 1102–1114. doi:10.1002/hbm.21497.
- Haixiang, G., Yijing, L., Shang, J., Mingyun, G., Yuanyue, H., Bing, G., 2017. Learning from class-imbalanced data: review of methods and applications. *Expert. Syst. Appl.* 73, 220–239. doi:10.1016/j.eswa.2016.12.035. <https://www.sciencedirect.com/science/article/pii/S0957417416307175>
- Hand, D.J., Till, R.J., 2001. A simple generalisation of the area under the roc curve for multiple class classification problems. *Mach. Learn.* 45 (2), 171–186. doi:10.1023/A:1010920819831.
- Hassabis, D., Kumaran, D., Summerfield, C., Botvinick, M., 2017. Neuroscience-inspired artificial intelligence. *Neuron* 95 (2), 245–258.
- Haxby, J.V., Gobbini, M.I., Furey, M.L., Ishai, A., Schouten, J.L., Pietrini, P., 2001. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science* 293 (5539), 2425–2430. doi:10.1126/science.1063736.
- He, H., Garcia, E.A., 2009. Learning from imbalanced data. *IEEE Trans. Knowl. Data Eng.* 21 (9), 1263–1284.
- Hebart, M.N., Gorgen, K., Haynes, J.-D., 2015. The decoding toolbox (tdt): a versatile software package for multivariate analyses of functional imaging data. *Front. Neuroinform.* 8, 88.
- Helmstaedter, M., 2015. The mutual inspirations of machine learning and neuroscience. *Neuron* 86 (1), 25–28.
- Hosseini, M., Powell, M., Collins, J., Callahan-Flintoft, C., Jones, W., Bowman, H., Wyble, B., 2020. I tried a bunch of things: the dangers of unexpected overfitting in classification of brain data. *Neurosci. Biobehav. Rev.* 119, 456–467. doi:10.1016/j.neubiorev.2020.09.036. <https://www.sciencedirect.com/science/article/pii/S0149763420305868>
- Hunter, J.D., 2007. Matplotlib: a 2d graphics environment. *Comput. Sci. Eng.* 9 (3), 90–95. doi:10.1109/MCSE.2007.55.
- Japkowicz, N., Stephen, S., 2002. The class imbalance problem: a systematic study. *Intell. Data Anal.* 6 (5), 429–449.
- Jeni, L.A., Cohn, J.F., De La Torre, F., 2013. Facing imbalanced data—recommendations for the use of performance metrics. In: 2013 Humaine association conference on affective computing and intelligent interaction. IEEE, pp. 245–251.
- Johnson, J.M., Khoshgoftaar, T.M., 2019. Survey on deep learning with class imbalance. *J. Big Data* 6 (1), 1–54.
- Kamalov, F., Denisov, D., 2020. Gamma distribution-based sampling for imbalanced data. *Knowl. Based Syst.* 207, 106368.
- Kelleher, J.D., Mac Namee, B., D'Arcy, A., 2015. Fundamentals of machine learning for predictive data analytics: algorithms. Worked examples, and case studies.
- Krawczyk, B., 2016. Learning from imbalanced data: open challenges and future directions. *Progr. Artif. Intell.* 5 (4), 221–232.
- Lee, K., Wu, X., Lee, Y., Lin, D.-T., Bhattacharyya, S.S., Chen, R., 2021. Neural decoding on imbalanced calcium imaging data with a network of support vector machines. *Adv. Rob.* 35 (7), 459–470. doi:10.1080/01691864.2020.1863259.
- Lemm, S., Blankertz, B., Dickhaus, T., Müller, K.-R., 2011. Introduction to machine learning for brain imaging. *Neuroimage* 56 (2), 387–399. doi:10.1016/j.neuroimage.2010.11.004. Multivariate Decoding and Brain Reading. <https://www.sciencedirect.com/science/article/pii/S1053811910014163>
- Liu, X.-Y., Wu, J., Zhou, Z.-H., 2008. Exploratory undersampling for class-imbalance learning. *IEEE Trans. Syst. Man Cybern., Part B (Cybern.)* 39 (2), 539–550.
- Luo, G., 2016. A review of automatic selection methods for machine learning algorithms and hyper-parameter values. *Netw. Model. Anal. Health Inform. Bioinform.* 5 (1), 18. doi:10.1007/s13721-016-0125-6.
- Luque, A., Carrasco, A., Martín, A., de Las Heras, A., 2019. The impact of class imbalance in classification performance metrics based on the binary confusion matrix. *Pattern Recognit.* 91, 216–231.
- Macpherson, T., Churchland, A., Sejnowski, T., DiCarlo, J., Kamitani, Y., Takahashi, H., Hikida, T., 2021. Natural and artificial intelligence: a brief introduction to the interplay between ai and neuroscience research. *Neural Netw.* 144, 603–613.
- McInnes, L., Healy, J., Melville, J., 2018. Umap: uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426.
- Myszczyńska, M.A., Ojames, P.N., Lacoste, A., Neil, D., Saffari, A., Mead, R., Hautbergue, G.M., Holbrook, J.D., Ferraiuolo, L., 2020. Applications of machine learning to diagnosis and treatment of neurodegenerative diseases. *Nature Rev. Neurol.* 16 (8), 440–456.
- Ojala, M., Garriga, G.C., 2010. Permutation tests for studying classifier performance. *J. Mach. Learn. Res.* 11 (6).
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E., 2011. Scikit-learn: machine learning in python. *J. Mach. Learn. Res.* 12, 2825–2830.
- Pereira, F., Mitchell, T., Botvinick, M., 2009. Machine learning classifiers and fmri: a tutorial overview. *Neuroimage* 45 (1), S199–S209.
- Prati, R.C., Batista, G.E., Silva, D.F., 2015. Class imbalance revisited: a new experimental setup to assess the performance of treatment methods. *Knowl. Inf. Syst.* 45 (1), 247–270. doi:10.1007/s10115-014-0794-3.
- Richards, B.A., Lillicrap, T.P., Beaudoin, P., Bengio, Y., Bogacz, R., Christensen, A., Clopath, C., Costa, R.P., de Berker, A., Ganguli, S., et al., 2019. A deep learning framework for neuroscience. *Nat. Neurosci.* 22 (11), 1761–1770.
- 3.2. Tuning the hyper-parameters of an estimator. https://scikit-learn.org/stable/modules/grid_search.html.
- Sagi, O., Rokach, L., 2018. Ensemble learning: a survey. *WIREs Data Min. Knowl. Discov.* 8 (4), e1249. doi:10.1002/widm.1249.
- Schalk, G., McFarland, D.J., Hinterberger, T., Birbaumer, N., Wolpaw, J.R., 2004. Bci2000: a general-purpose brain-computer interface (bci) system. *IEEE Trans. Biomed. Eng.* 51 (6), 1034–1043.
- Shafra, M.A., Tyler, L.K., Dixon, M., Taylor, J.R., Rowe, J.B., Cusack, R., Calder, A.J., Marslen-Wilson, W.D., Duncan, J., Dalgleish, T., et al., 2014. The cambridge centre for ageing and neuroscience (cam-can) study protocol: a cross-sectional, lifespan, multidisciplinary examination of healthy cognitive ageing. *BMC Neurol.* 14 (1), 1–25.
- Siblini, W., Fréry, J., He-Guelton, L., Oblé, F., Wang, Y.-Q., 2020. Master your metrics with calibration. In: International Symposium on Intelligent Data Analysis. Springer, pp. 457–469.
- Somasundaram, A., Reddy, U.S., 2017. Modelling a stable classifier for handling large scale data with noise and imbalance. In: 2017 International Conference on Computational Intelligence in Data Science (ICCIDS). IEEE, pp. 1–6.
- Straube, S., Krell, M.M., 2014. How to evaluate an agent's behavior to infrequent events? reliable performance estimation insensitive to class distribution. *Front. Comput. Neurosci.* 8, 43.
- Sun, Y., Wong, A.K., Kamel, M.S., 2009. Classification of imbalanced data: a review. *Int. J. Pattern Recognit. Artif. Intell.* 23 (04), 687–719.
- Tan, T.Z., Ng, G.S., Quek, C., 2007. Complementary learning fuzzy neural network: an approach to imbalanced dataset. In: 2007 International Joint Conference on Neural Networks. IEEE, pp. 2306–2311.
- Taylor, J.R., Williams, N., Cusack, R., Auer, T., Shafra, M.A., Dixon, M., Tyler, L.K., Henson, R.N., et al., 2017. The cambridge centre for ageing and neuroscience (cam-can) data repository: structural and functional mri, meg, and cognitive data from a cross-sectional adult lifespan sample. *Neuroimage* 144, 262–269.
- Thabtah, F., Hammoud, S., Kamalov, F., Gonsalves, A., 2020. Data imbalance in classification: experimental evaluation. *Inf. Sci. (Nij)* 513, 429–441.
- Van De Ruit, M., Grey, M., 2019. The large type 1 error associated with responder analyses. *Brain Stimul.* 12 (2), 525–526. doi:10.1016/j.brs.2018.12.729.
- Varoquaux, G., Colliot, O., 2022. Evaluating Machine Learning Models and Their Diagnostic Value. Machine Learning for Brain Disorders. <https://hal.archives-ouvertes.fr/hal-03682454>
- Wang, Q., Luo, Z., Huang, J., Feng, Y., Liu, Z., 2017. A novel ensemble method for imbalanced data learning: bagging of extrapolation-smote svm. *Comput. Intell. Neurosci.* 2017.
- Wardhani, N.W.S., Rochayani, M.Y., Iriany, A., Sulistyono, A.D., Lestantyo, P., 2019. Cross-validation metrics for evaluating classification performance on imbalanced data. In: 2019 international conference on computer, control, informatics and its applications (IC3INA). IEEE, pp. 14–18.
- Waskom, M.L., 2021. Seaborn: statistical data visualization. *J. Open Source Softw.* 6 (60), 3021. doi:10.21105/joss.03021.
- Yang, G.R., Wang, X.-J., 2020. Artificial neural networks for neuroscientists: a primer. *Neuron* 107 (6), 1048–1070.
- Yang, L., Shami, A., 2020. On hyperparameter optimization of machine learning algorithms: theory and practice. *Neurocomputing* 415, 295–316. doi:10.1016/j.neucom.2020.07.061. <https://www.sciencedirect.com/science/article/pii/S09525231220311693>
- Zhu, M., Xia, J., Jin, X., Yan, M., Cai, G., Yan, J., Ning, G., 2018. Class weights random forest algorithm for processing class imbalanced medical data. *IEEE Access* 6, 4641–4652.