

# UC Berkeley

## UC Berkeley Electronic Theses and Dissertations

### Title

Robust, automated sleep scoring by a compact neural network with distributional shift correction

### Permalink

<https://escholarship.org/uc/item/2jt7v8pd>

### Author

Barger, Zeke

### Publication Date

2020

Peer reviewed|Thesis/dissertation

Robust, automated sleep scoring by a compact neural network  
with distributional shift correction

by

Zeke K. Barger

A dissertation submitted in partial satisfaction of the

requirements for the degree of

Doctor of Philosophy

in

Neuroscience

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor Yang Dan, Chair

Professor Lance Kriegsfeld

Professor Ehud Isacoff

Assistant Professor Stephan Lammel

Spring 2020

## Abstract

Robust, automated sleep scoring by a compact neural network  
with distributional shift correction

by

Zeke K. Barger

Doctor of Philosophy in Neuroscience

University of California, Berkeley

Professor Yang Dan, Chair

Accurately determining the sleep stage of experimental subjects is a key step in sleep research. Despite years of research into automated methods for scoring rodent sleep recordings, most scoring is still performed manually. Here, I present an automated, machine learning-based sleep scoring method that avoids the subjective and labor-intensive task of manual scoring. In the first chapter, I review recent advances in the field of sleep scoring. New algorithms have, over time, extracted more and more useful information from underlying physiological signals used as inputs. However, inter-laboratory and inter-subject differences have thus far prevented any single automated method from being widely applicable.

In the second chapter, I present a feature scaling algorithm, mixture  $z$ -scoring, that can eliminate many of these differences. Importantly, it also preserves changes in the amount of time a given subject spends in each sleep stage, which is not attainable using existing algorithms. I then present a neural network architecture which efficiently learns to score sleep from spectrograms of electroencephalogram recordings and evaluate it using a large, high-quality dataset. When mixture  $z$ -scoring is used as a preprocessing step, the network achieves state-of-the-art performance. I also introduce a free, open-source software package that allows even novice users to make use of the network and mixture  $z$ -scoring. This work is presented in the form of a published, first-author manuscript.

In the final chapter, I discuss the limitations of the scoring algorithm and its potential application for scoring data from other species. I also examine some remaining challenges in the field of sleep scoring as well as their possible solutions. As a whole, this work provides computational tools that are designed to meet the data processing needs of the sleep research community.

*For Stuart Curtis Barger*

# Contents

<b>Contents</b>	<b>ii</b>
<b>List of Figures</b>	<b>iii</b>
<b>Acknowledgements</b>	<b>iv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Sleep scoring . . . . .	2
1.3 Distributional shift . . . . .	7
References . . . . .	9
<b>2 Robust, automated sleep scoring by a compact neural network with distributional shift correction</b>	<b>13</b>
2.1 Abstract . . . . .	13
2.2 Introduction . . . . .	14
2.3 Results . . . . .	15
2.3.1 Mixture $z$ -scoring corrects for distributional shift in simulated data .	15
2.3.2 Mixture $z$ -scoring reduces bias when classifying mouse <i>in vivo</i> data .	18
2.3.3 Validation of SS-ANN . . . . .	19
2.3.4 AccuSleep: free, open-source software for automated sleep scoring . .	20
2.4 Discussion . . . . .	21
2.5 Methods . . . . .	23
2.5.1 Polysomnographic recordings . . . . .	23
2.5.2 Sleep scoring algorithm . . . . .	25
References . . . . .	28
<b>3 Conclusion</b>	<b>42</b>
3.1 Closing remarks . . . . .	42
3.2 Outlook . . . . .	43
References . . . . .	43

# List of Figures

1.1	Usage of sleep scoring methods . . . . .	2
1.2	Classification with a decision tree . . . . .	3
1.3	Classification with a support vector machine . . . . .	5
1.4	Convolutional neural network architecture . . . . .	6
1.5	Distributional shift and $z$ -scoring . . . . .	8
2.1	Overview of the signal collection process for sleep scoring in mice . . . . .	30
2.2	Correcting for distributional shift prevents a false negative in a simple model . .	32
2.3	Comparison of sleep scoring algorithms . . . . .	34
2.4	Validation of SS-ANN . . . . .	36
2.5	AccuSleep interface for manual sleep scoring . . . . .	38
2.6	AccuSleep interface for automated sleep scoring . . . . .	40

## Acknowledgments

I would first like to thank my advisor, Yang Dan, for the mentorship, guidance, and encouragement that made this work possible. It has been a pleasure working with the members of the Dan lab: Franz Weber and Johnny Do, who helped me get started in the lab when I knew nothing about sleep research; Chenyan Ma, Yuanyuan Yao, Danqian Liu, and Peng Zhong, who provided much-needed advice on practical matters; and Zhe Zhang, who was always available for thoughtful discussions. I truly enjoyed my collaboration with Charles Frye, who helped me develop some of the key ideas presented in this work. I have been lucky to meet so many amazing people during my time in graduate school, and this would hardly have been possible without them either: I thank Gloria Yu, Krisha Aghi, the boys at the wall, the members of Stumphole, and everyone in the HWNI community. And most of all, I thank my family for their love.

# Chapter 1

## Introduction

### 1.1 Background

Sleep is critical for our mental and physical well-being and plays a role in immune function [1], cognition [2], and the processing of emotional experiences [3, 4]. Disturbances in sleep are associated with a wide range of mental health disorders [5]. Furthermore, a number of key scientific questions about sleep remain unresolved. The fact that sleep, or a sleep-like state, is observed in all vertebrates and possibly all animals [6], together with the fact that organisms are vulnerable during sleep, has motivated a search for features of sleep that make it adaptive or necessary. Possibilities include clearing cellular waste [7], regulation of synaptic strength [8], and memory consolidation [9]. The neural mechanisms of sleep and its homeostatic regulation are also active areas of research [10, 11].

Addressing these clinical and scientific questions requires methods that can continuously monitor the brain state of a patient or subject. The outward signs of sleep are reduced sensitivity to stimuli, stereotypical posture, and few movements [12]. The invention of the electroencephalogram (EEG) by Hans Berger in 1924 allowed the first window onto brain activity during this state. With the introduction of polysomnography (PSG)—the simultaneous recording of EEG and other signals such as the electromyogram (EMG) and electrooculogram (EOG)—it soon became clear that not only did sleep have a distinct pattern of brain electrical activity from wakefulness, but also that sleep could be broken down into several stages that were consistent across subjects and nights. The first standardized guidelines for scoring these stages were published in 1968 [13].

Today, the American Academy of Sleep Medicine maintains a comprehensive manual for scoring human sleep [14]. The rules for scoring each stage are complex but can be summarized as follows: Stage N1 marks the transition into sleep, and is characterized by drowsiness, slow eye movements, and theta (4-7 Hz) EEG activity. Stage N2, or light sleep, is marked by phasic EEG phenomena called spindles and K-complexes. Stage N3, or deep



sleep, is characterized by high-amplitude delta (0.5-2 Hz) EEG activity. Rapid eye movement (REM) sleep features bursts of eye movements, very low muscle tone, and wake-like EEG activity.

Model organisms used for research show patterns of sleep with varying levels of similarity to humans, and rules for scoring sleep in each organism are determined by a consensus among researchers. Mice have been used extensively in sleep research, in part due to the development of genetic tools that allow for targeted manipulation and interrogation of neuronal populations [15]. Only two sleep stages are recognized in mice: REM and non-REM (see section 2.2). However, reliably detecting these stages and distinguishing them from wakefulness is not trivial and is, itself, an active field of research. In this dissertation, I review the progress and challenges in this field and present a novel method for mouse sleep scoring.

## 1.2 Sleep scoring

Manual inspection of PSG data by expert scorers remains the most commonly used method for mouse sleep scoring (Fig 1.1). The low throughput of this approach has motivated many years of research into automated scoring methods, which promise not only to increase classification speed, but also to reduce experimenter bias and eliminate variability due to inter-scorer differences. Such algorithms must be highly accurate if they are to be at all useful because common measurements, such as the mean duration of bouts of each stage, are sensitive to noisy predictions. They must also operate at high temporal resolution in order to measure the effect of rapid experimental manipulations, such as optogenetic stimuli.

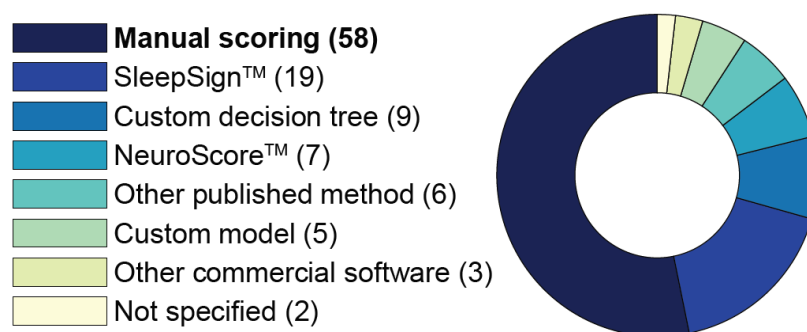


Figure 1.1: **Usage of sleep scoring methods.** Each segment represents the number of papers published between 2017 and 2019 that reported using a given sleep scoring method. The word 'model' indicates a machine learning-based classifier. Note that the SleepSign algorithm is essentially a decision tree. Source: PubMed search for 'mouse wake sleep rem'.

Decision trees [16–20] are the second-most common approach to sleep scoring. In this class of models, the input data are successively split based on thresholds applied to individual features (Fig 1.2A). Here, the word “feature” refers to something that is measured about each observation, such as the age of a subject or a pixel in an image. Sleep scoring methods that rely on decision trees typically have a fixed tree structure, but adjust the thresholds for each subject or recording. Some advantages of these models are that they are perfectly interpretable—i.e., one can follow the tree to understand each prediction—and can classify new observations faster than manual scoring. However, manually adjusting the parameters is not only time-consuming, but also introduces subjectivity into the analysis. Additionally, shallow decision trees may lack the complexity to construct an adequate decision boundary between different classes of observations (Fig 1.2B).

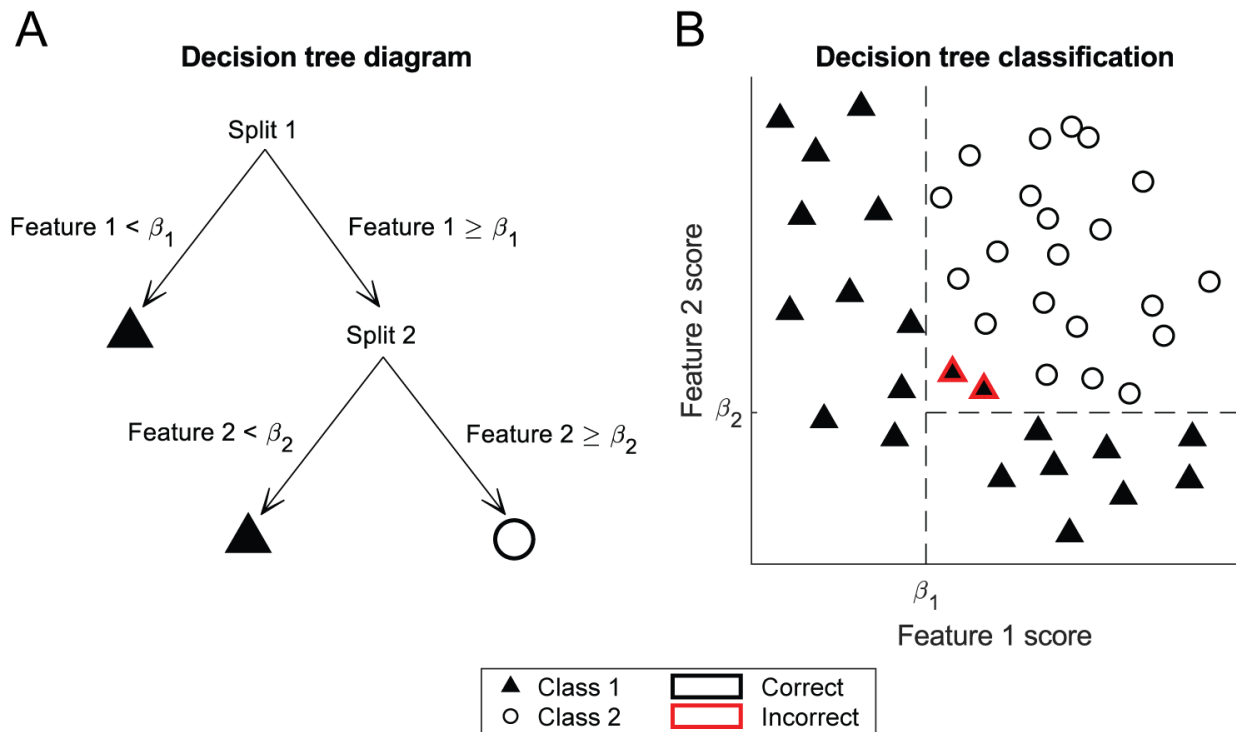


Figure 1.2: **Classification with a decision tree.** A: Example tree diagram. At each internal node, observations are split into two groups based on whether a single feature of the observations is above or below a threshold,  $\beta$ . Observations in a leaf node are assigned a label and are not split any further. In this tree, the first split assigns observations where feature 1 is less than  $\beta_1$  to class 1. The remaining observations are split again so that observations where feature 2 is less than  $\beta_2$  are assigned to class 1 and the remainder are assigned to class 2. B: Application of the decision tree classifier in A to a toy dataset. Dashed lines represent the thresholds in the decision tree. Each marker represents an observation, with class denoted by the marker shape. Two observations from class 1 are incorrectly classified by this tree.

Machine learning techniques have the ability to overcome these limitations. In supervised machine learning for classification, a learning algorithm uses a set of labeled examples (a “training set”) to automatically adjust the parameters of a model so that unlabeled inputs (a “test set”) can be correctly classified. The learning algorithm therefore allows the complexity of the model and the dimensionality (i.e., the number of features) of the input data to be very high, in contrast to the decision trees described earlier.

Support-vector machines (SVMs) are one class of supervised learning models that have been applied to rodent sleep scoring [21, 22]. SVM algorithms find the decision boundary that separates two classes of observations by the largest possible margin. To account for the case where the classes are not fully separable (Fig 1.3A), a hyperparameter,  $C$ , is used to control the trade-off between reducing the number of errors and maximizing the size of the margin. One advantage of SVMs over decision trees is that SVMs can find nonlinear decision boundaries with the appropriate setting of another hyperparameter, called the kernel, which maps observations into a different feature space (Fig 1.3B). However, a drawback of SVMs is that unless the value of  $C$  is chosen carefully, the performance of SVMs suffers when the feature set is large and includes many features that do not contribute useful information. This means that the most informative features should be identified before training the model, and for sleep scoring, it is not always clear which features of the EEG signal are most informative.

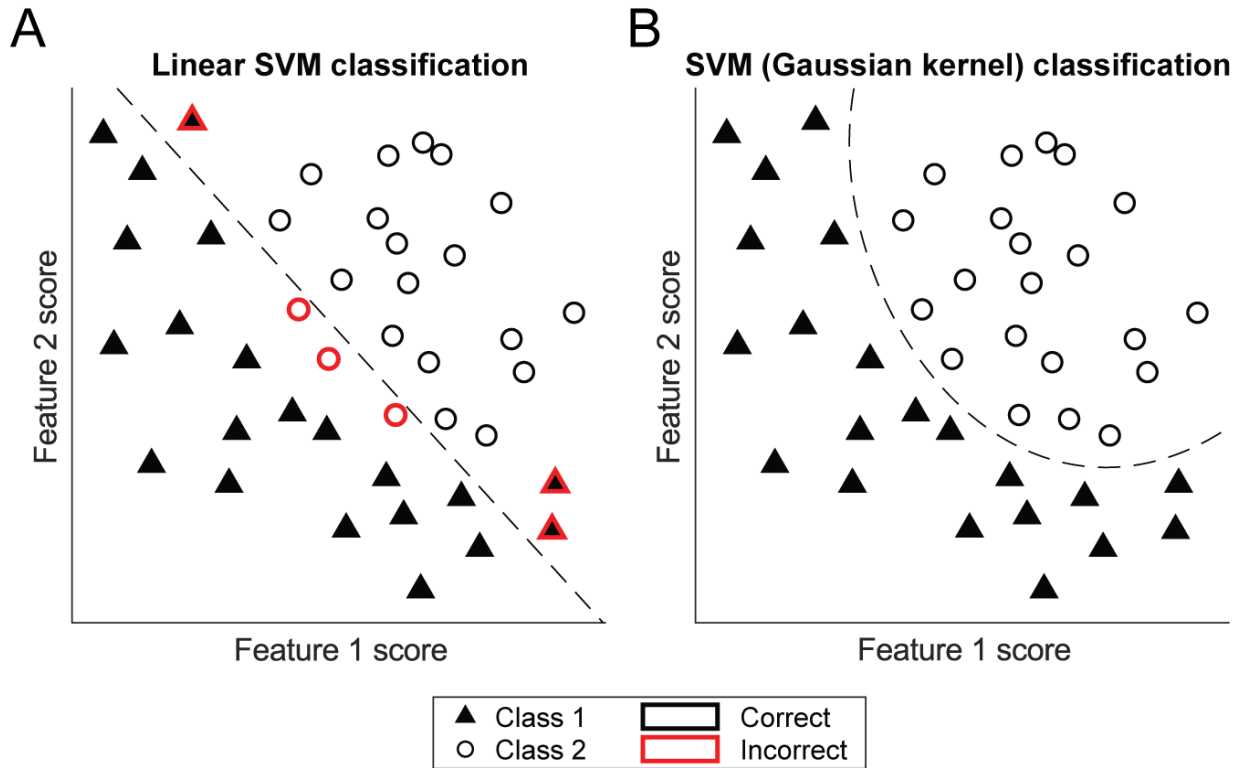


Figure 1.3: **Classification with a support vector machine.** A: Application of a linear support vector machine (SVM) to the toy dataset. Because the data are not linearly separable, several observations from each class lie on the wrong side of the decision boundary and are incorrectly classified. B: Application of a SVM with a Gaussian kernel to the toy dataset. Thanks to the nonlinear decision boundary, all observations are correctly classified.

Spectrograms, which represent the frequency content of the EEG signal at each moment in time (Fig 2.1), may be an effective way to capture many useful features simultaneously. Several recent studies have explored the use of convolutional neural networks (CNNs) to score PSG data in the form of spectrograms [23, 24]. CNNs are artificial neural networks that efficiently learn the informative features of images and are therefore widely used for image classification. The key element of a CNN is the convolutional layer, which calculates the similarity between a filter and every patch of an image and stores the result in a feature map (Fig 1.4). The feature maps are typically pooled (downsampled) to reduce the number of parameters in the model. Deeper convolutional layers take pooled feature maps from earlier layers as input, capturing progressively higher-level image features. Finally, a fully connected layer and softmax layer classify the images based on features detected by the convolutional layers. The weights in the fully connected layer and the filters in the convolutional layers are adjusted by a learning algorithm based on labeled images in the training set. Thus, by

automatically learning the informative features of spectrograms, CNNs can score PSG data without requiring the features to be selected in advance.

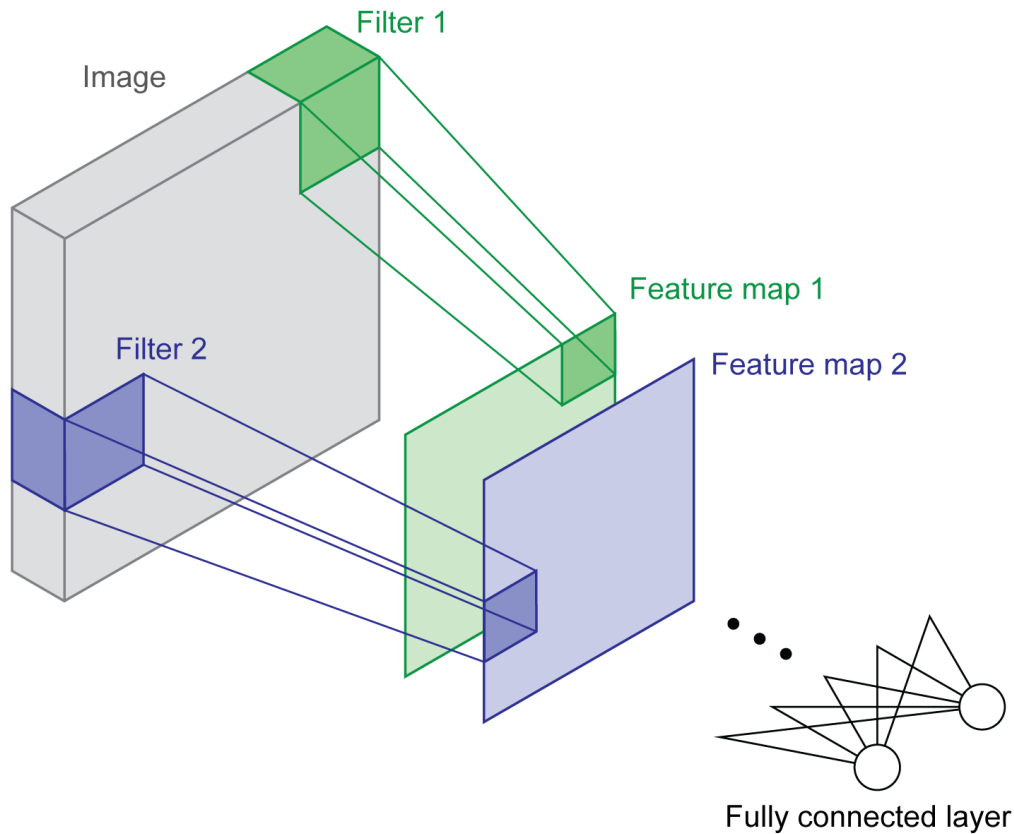


Figure 1.4: **Convolutional neural network (CNN) architecture.** A convolutional layer convolves one or more filters (green and blue volumes) with an image (gray volume) to produce feature maps (green and blue squares). The feature maps are typically subjected to batch normalization, a nonlinearity (such as the ReLU function), and pooling before being used as the input to a deeper convolutional layer. A fully connected layer learns combinations of the features identified by the convolutional layers to perform classification.

All of the methods described above offer significant speed improvements over manual scoring and are reported to achieve accuracy in the range of 89-96%. The machine learning approaches are particularly promising because they do not seem to require continual re-adjusting of their parameters. However, usage of these methods remains low (Fig 1.1). One possible explanation is that their ability to generalize, or perform well on new data, is poor. As I explain in the next section, variability in mouse PSG data indeed poses a challenge to automated sleep scoring algorithms.

### 1.3 Distributional shift

Differences between training and test sets, called dataset shift, are commonly encountered when applying machine learning algorithms [25], especially in biology when data are collected with different equipment in different laboratories [26, 27]. Distributional shift (or domain shift [28]) is a type of dataset shift where the values of some or all data features are shifted in some way. To illustrate why this is problematic, imagine that the data points in Fig 1.2B represent a training set and the dashed line represents the fitted decision boundary of a classification model. If a test set were exactly like the training set except shifted rightward, many of the data points belonging to class 1 would be on the wrong side of the decision boundary.

Approaches for counteracting distributional shift belong to the field of transfer learning. One application of transfer learning is to adapt a model trained in one domain so that it can perform well in another domain. For example, the model could be trained to detect cats in photographs taken indoors, and would require modifications to detect cats in photographs taken outdoors. The simplest solution is, of course, to train a new model for each domain. However, mouse PSG recordings show variability not just between laboratories (due to differences in recording equipment), but also on an individual level. The content of EEG and EMG signals can vary from mouse to mouse due to the position or impedance of the recording electrodes, or even due to mouse genotype [29]. This means that each mouse is, essentially, a new domain, and labeling a sufficiently large training set for each mouse defeats the purpose of training individualized models.

Another way to counteract distributional shift is to find a representation of the data that is invariant across domains [30]. In some cases, this is achievable by normalizing the feature values from each domain—i.e., finding a common scale for each feature. Several recent publications have proposed that normalization could be an effective way to address the variability in mouse PSG data [24, 31]. Specifically, Miladinović et al. applied  $z$ -scoring (also called standardization) to every feature on a per-recording basis.

$z$ -scores are calculated by subtracting the sample mean and dividing by the sample standard deviation (Eq 2.1). To visualize how  $z$ -scoring transforms a dataset, consider the hypothetical experimental results in Fig 1.5. Say that we collect observations from three subjects, and that the observations fall into three classes which are associated with different values of the feature we are measuring. For subjects 1 and 2, the proportion of observations from each class is the same, but the measurements are affected by an affine transformation (i.e., one composed of shifting and scaling)—a case of domain shift. After  $z$ -scoring the measurements (panels B and D), the data from the two subjects are indeed brought to a common scale.

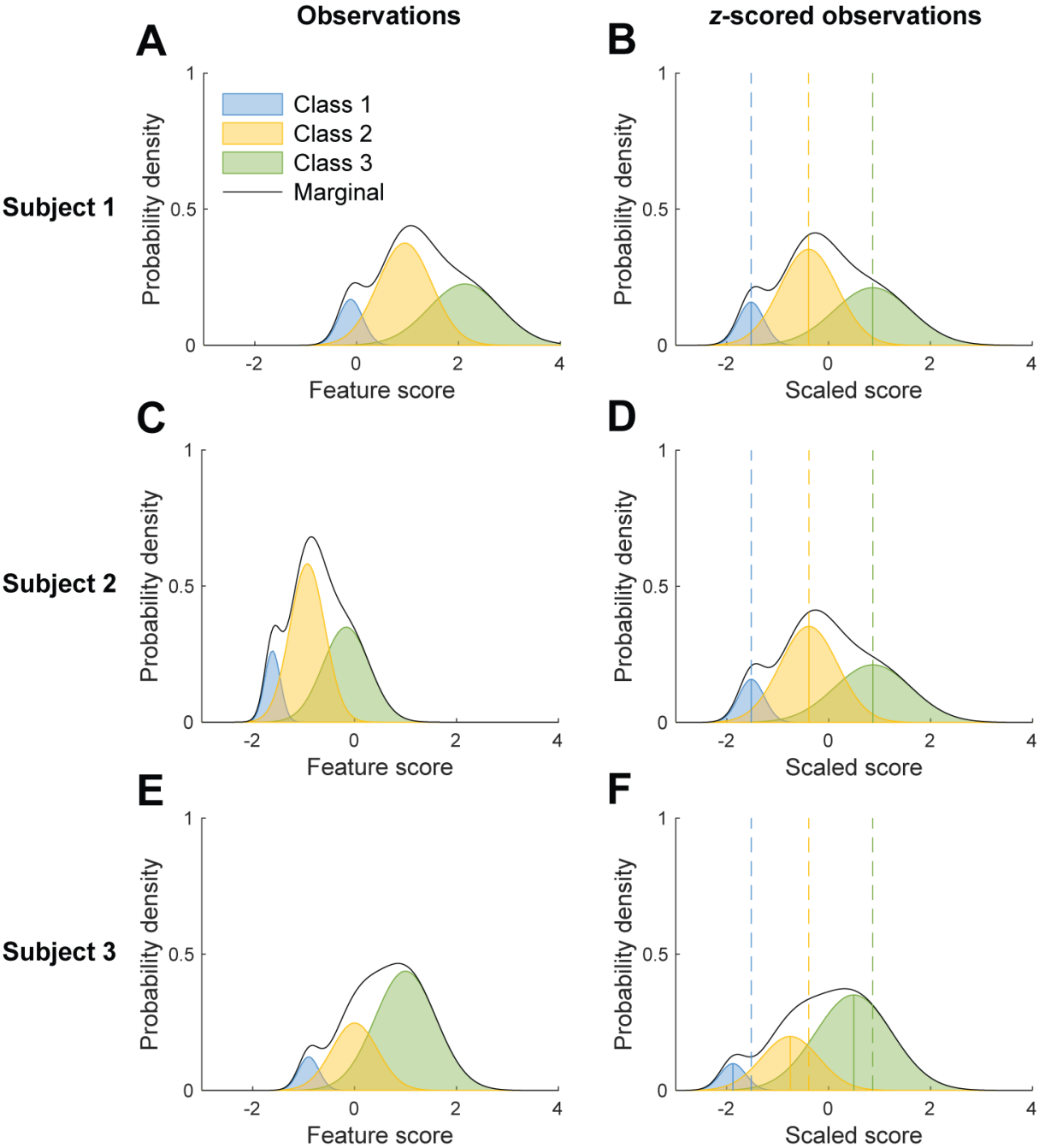


Figure 1.5: **Distributional shift and  $z$ -scoring.** This figure uses a hypothetical experiment to demonstrate how  $z$ -scoring affects distributions of feature scores under different kinds of distributional shift. Each row represents data from one subject. The first column contains probability density plots of the original observations. The second column contains probability density plots of the observations following  $z$ -scoring. (Continued on next page)

Figure 1.5: (Continued from previous page) Dashed lines mark the class-conditional means following  $z$ -scoring for subject 1. The distributions of observations from subjects 1 and 2 differ only by an affine transformation (an example of domain shift), while subject 3 also shows prior probability shift. As a result, only data from subjects 1 and 2 are brought to a common scale by  $z$ -scoring.

However, PSG data are subject to another form of dataset shift that cannot be corrected by  $z$ -scoring: prior probability shift [32]. This occurs when the proportion of observations belonging to each class differs between the training and test sets. For example, one might train a classifier on a dataset that contains equal quantities of wakefulness and sleep, and then apply it to a recording from a sleep deprivation experiment. Fig 1.5E shows data from a third hypothetical subject with different proportions of observations from the first two. In panel F, it is clear that  $z$ -scoring observations subject to this prior probability shift did not bring the data onto a common scale as desired. The situation where multiple forms of dataset shift are present—in this case, domain shift and prior probability shift—is recognized as the most difficult to solve [33]. In summary, dataset shift poses a challenge to machine learning algorithms for sleep scoring that has not yet been adequately addressed.

In the next chapter, I use both simulated and real mouse PSG data to demonstrate how dataset shift is problematic for contemporary automated scoring methods. I then present an algorithm that brings features of PSG data into a common scale, thus allowing a convolutional neural network to achieve state-of-the-art classification accuracy even in the presence of dataset shift.

## References

- [1] Besedovsky, L., Lange, T., & Born, J. (2012). Sleep and immune function. *Pflügers Archiv-European Journal of Physiology*, *463*(1), 121–137.
- [2] Walker, M. P. (2009). The role of sleep in cognition and emotion. *Annals of the New York Academy of Sciences*, *1156*(1), 168–197.
- [3] van der Helm, E., Yao, J., Dutt, S., Rao, V., Saletin, J. M., & Walker, M. P. (2011). REM sleep depotentiates amygdala activity to previous emotional experiences. *Current Biology*, *21*(23), 2029–2032. <https://doi.org/10.1016/j.cub.2011.10.052>
- [4] Wassing, R., Lakbila-Kamal, O., Ramautar, J. R., Stoffers, D., Schalkwijk, F., & Someren, E. J. [ (2019). Restless REM sleep impedes overnight amygdala adaptation. *Current Biology*, *29*(14), 2351–2358.e4. <https://doi.org/10.1016/j.cub.2019.06.034>
- [5] Anderson, K. N., & Bradley, A. J. (2013). Sleep disturbance in mental health problems and neurodegenerative disease. *Nature and Science of Sleep*, *5*, 61.



- [6] Nath, R. D., Bedbrook, C. N., Abrams, M. J., Basinger, T., Bois, J. S., Prober, D. A., Sternberg, P. W., Gradinaru, V., & Goentoro, L. (2017). The jellyfish *cassiopea* exhibits a sleep-like state. *Current Biology*, *27*(19), 2984–2990.
- [7] Xie, L., Kang, H., Xu, Q., Chen, M. J., Liao, Y., Thiyagarajan, M., O’Donnell, J., Christensen, D. J., Nicholson, C., Iliff, J. J., Takano, T., Deane, R., & Nedergaard, M. (2013). Sleep drives metabolite clearance from the adult brain. *Science*, *342*(6156), 373–377. <https://doi.org/10.1126/science.1241224>
- [8] Tononi, G., & Cirelli, C. (2014). Sleep and the price of plasticity: From synaptic and cellular homeostasis to memory consolidation and integration. *Neuron*, *81*(1), 12–34. <https://doi.org/10.1016/j.neuron.2013.12.025>
- [9] Stickgold, R. (2005). Sleep-dependent memory consolidation. *Nature*, *437*(7063), 1272–1278.
- [10] Faraguna, U., Vyazovskiy, V. V., Nelson, A. B., Tononi, G., & Cirelli, C. (2008). A causal role for brain-derived neurotrophic factor in the homeostatic regulation of sleep. *Journal of Neuroscience*, *28*(15), 4088–4095. <https://doi.org/10.1523/JNEUROSCI.5510-07.2008>
- [11] Liu, D., & Dan, Y. (2019). A motor theory of sleep-wake control: Arousal-action circuit. *Annual Review of Neuroscience*, *42*(1), 27–46. <https://doi.org/10.1146/annurev-neuro-080317-061813>
- [12] Harrington, J., & Lee-Chiong, T. (2012). Basic biology of sleep [Sleep Medicine and Dentistry]. *Dental Clinics of North America*, *56*(2), 319–330. <https://doi.org/10.1016/j.cden.2012.01.005>
- [13] Deak, M., & Epstein, L. J. (2009). The history of polysomnography. *Sleep Medicine Clinics*, *4*(3), 313–321.
- [14] Berry, R. B., Brooks, R., Gamaldo, C. E., Harding, S. M., Marcus, C., Vaughn, B. V. Et al. (2012). The AASM manual for the scoring of sleep and associated events. *Rules, Terminology and Technical Specifications, Darien, Illinois, American Academy of Sleep Medicine*, *176*, 2012.
- [15] Ma, C., Zhong, P., Liu, D., Barger, Z. K., Zhou, L., Chang, W.-C., Kim, B., & Dan, Y. (2019). Sleep regulation by neurotensinergic neurons in a thalamo-amygdala circuit. *Neuron*, *103*(2), 323–334.
- [16] Stephenson, R., Caron, A. M., Cassel, D. B., & Kostela, J. C. (2009). Automated analysis of sleep-wake state in rats. *Journal of Neuroscience Methods*, *184*(2), 263–274. <https://doi.org/10.1016/j.jneumeth.2009.08.014>
- [17] Kohtoh, S., Taguchi, Y., Matsumoto, N., Wada, M., Huang, Z., & Urade, Y. (2008). Algorithm for sleep scoring in experimental animals based on fast fourier transform power spectrum analysis of the electroencephalogram. *Sleep and Biological Rhythms*, *6*(3), 163–171. <https://doi.org/10.1111/j.1479-8425.2008.00355.x>
- [18] Bastianini, S., Berteotti, C., Gabrielli, A., Vecchio, F. D., Amici, R., Alexandre, C., Scammell, T. E., Gazea, M., Kimura, M., Martire, V. L., Silvani, A., & Zoccoli, G. (2014). SCOPRISM: A new algorithm for automatic sleep scoring in mice. *Journal of Neuroscience Methods*, *235*, 277–284. <https://doi.org/10.1016/j.jneumeth.2014.07.018>

- [19] Kreuzer, M., Polta, S., Gapp, J., Schuler, C., Kochs, E., & Fenzl, T. (2015). Sleep scoring made easy—semi-automated sleep analysis software and manual rescoring tools for basic sleep research in mice. *MethodsX*, 2, 232–240. <https://doi.org/10.1016/j.mex.2015.04.005>
- [20] Gross, B. A., Walsh, C. M., Turakhia, A. A., Booth, V., Mashour, G. A., & Poe, G. R. (2009). Open-source logic-based automated sleep scoring software using electrophysiological recordings in rats. *Journal of Neuroscience Methods*, 184(1), 10–18. <https://doi.org/10.1016/j.jneumeth.2009.07.009>
- [21] Shantilal, Donohue, K. D., & O’Hara, B. F. (2008). SVM for automatic rodent sleep-wake classification, In *IEEE SoutheastCon 2008*.
- [22] Crisler, S., Morrissey, M. J., Anch, A. M., & Barnett, D. W. (2008). Sleep-stage scoring in the rat using a support vector machine. *Journal of Neuroscience Methods*, 168(2), 524–534. <https://doi.org/10.1016/j.jneumeth.2007.10.027>
- [23] Vilamala, A., Madsen, K. H., & Hansen, L. K. (2017). Deep convolutional neural networks for interpretable analysis of EEG sleep stage scoring. *arXiv e-prints*, arXiv 1710.00633.
- [24] Miladinović, D., Muheim, C., Bauer, S., Spinnler, A., Noain, D., Bandarabadi, M., Gallusser, B., Krummenacher, G., Baumann, C., Adamantidis, A., Brown, S. A., & Buhmann, J. M. (2019). SPINDLE: End-to-end learning from EEG/EMG to extrapolate animal sleep scoring across experimental settings, labs and species. *PLOS Computational Biology*, 15(4), 1–30. <https://doi.org/10.1371/journal.pcbi.1006968>
- [25] Yin, D., Lopes, R. G., Shlens, J., Cubuk, E. D., & Gilmer, J. (2019). A Fourier perspective on model robustness in computer vision, In *Advances in neural information processing systems*.
- [26] Borgwardt, K. M., Gretton, A., Rasch, M. J., Kriegel, H.-P., Schölkopf, B., & Smola, A. J. (2006). Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14), e49–e57.
- [27] Stacke, K., Eilertsen, G., Unger, J., & Lundström, C. (2019). A closer look at domain shift for deep learning in histopathology. *arXiv preprint arXiv:1909.11575*.
- [28] Quionero-Candela, J., Sugiyama, M., Schwaighofer, A., & Lawrence, N. D. (2009). *Dataset shift in machine learning*. The MIT Press.
- [29] Franken, P., Malafosse, A., & Tafti, M. (1998). Genetic variation in eeg activity during sleep in inbred mice. *American Journal of Physiology-Regulatory, Integrative and Comparative Physiology*, 275(4), R1127–R1137.
- [30] Blitzer, J., McDonald, R., & Pereira, F. (2006). Domain adaptation with structural correspondence learning, In *Proceedings of the 2006 conference on empirical methods in natural language processing*.
- [31] Katsageorgiou, V., Lassi, G., Tucci, V., Murino, V., & Sona, D. (2015). Sleep-stage scoring in mice: The influence of data pre-processing on a system’s performance, In *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. <https://doi.org/10.1109/EMBC.2015.7318433>

- [32] Moreno-Torres, J. G., Raeder, T., Alaiz-Rodriguez, R., Chawla, N. V., & Herrera, F. (2012). A unifying view on dataset shift in classification. *Pattern recognition*, 45(1), 521–530.
- [33] Kouw, W. M., & Loog, M. (2018). An introduction to domain adaptation and transfer learning. *arXiv preprint arXiv:1812.11806*.

## Chapter 2

# Robust, automated sleep scoring by a compact neural network with distributional shift correction

**Zeke Barger, Charles G. Frye, Danqian Liu, Yang Dan, and Kristofer E. Bouchard**

This chapter, in full, is a replication of the material as it appears in Barger Z, Frye CG, Liu D, Dan Y, Bouchard KE (2019) Robust, automated sleep scoring by a compact neural network with distributional shift correction. PLoS ONE 14(12): e0224642.

### 2.1 Abstract

Studying the biology of sleep requires the accurate assessment of the state of experimental subjects, and manual analysis of relevant data is a major bottleneck. Recently, deep learning applied to electroencephalogram and electromyogram data has shown great promise as a sleep scoring method, approaching the limits of inter-rater reliability. As with any machine learning algorithm, the inputs to a sleep scoring classifier are typically standardized in order to remove distributional shift caused by variability in the signal collection process. However, in scientific data, experimental manipulations introduce variability that should not be removed. For example, in sleep scoring, the fraction of time spent in each arousal state can vary between control and experimental subjects. We introduce a standardization method, *mixture z-scoring*, that preserves this crucial form of distributional shift. Using both a simulated experiment and mouse *in vivo* data, we demonstrate that a common standardization method used by state-of-the-art sleep scoring algorithms introduces systematic bias, but that mixture *z-scoring* does not. We present a free, open-source user interface that uses a compact neural network and mixture *z-scoring* to allow for rapid sleep scoring with accuracy that compares well to contemporary methods. This work provides a set of computational tools for the robust automation of sleep scoring.

## 2.2 Introduction

Sleep is a fundamental animal behavior and has long been the subject of intensive basic and clinical research. Mice are commonly chosen as a model organism for sleep research thanks to the wide range of genetic tools that enable manipulation of sleep-relevant neuronal ensembles and characterization of sleep phenotypes. In order to measure the effect of an experimental manipulation on the quantity or timing of sleep stages, it is necessary to score a subject’s sleep stage at each point in time. In mice, each epoch is typically assigned to one of three stages based on patterns of electroencephalogram (EEG) and electromyogram (EMG) activity (Fig 2.1): rapid eye movement (REM) sleep, with a high ratio of theta (6-8 Hz) to delta (1-4 Hz) EEG activity and low muscle tone; non-REM (NREM) sleep, with a low theta/delta ratio and low muscle tone; and wakefulness, with high muscle tone and high-frequency, low-amplitude EEG activity.

Manual inspection of the EEG and EMG signals remains the most widely used method for mouse sleep scoring. However, this process is time-intensive and therefore scales poorly with the number of subjects and recordings, motivating efforts to develop scoring methods that are partially or completely automated. Shallow decision trees [1–5], which require a user to define thresholds in a low-dimensional feature space (e.g., theta/delta ratio and EMG activity), are one such approach. However, since the three classes are not entirely separable in these low-dimensional spaces, efforts have been made to build classifiers that use machine learning to exploit a larger number of hand-tuned features [6–8]. Most recently, there have been successes in using models trained directly on EEG/EMG data without feature engineering, either in the form of spectrograms [9, 10] or unprocessed signals [11, 12]. The accuracy of these methods on held-out test sets can be close to the inter-rater reliability of expert scorers [10], suggesting that further feature or architecture engineering of EEG/EMG-based sleep scoring algorithms will yield diminishing returns.

Generalization to test datasets that differ from training datasets, however, remains a concern, both for machine learning in general [13] and for automated sleep scoring in particular. Changes in the distribution of test data, called *distributional shift*, can cause misclassification errors when items of one class in the test set artifactually resemble, due to the shift, items of another class in the training set. Simple forms of distributional shift can be solved by a standardization procedure, such as  $z$ -scoring. In a review of automated sleep scoring methods, Katsageorgiou et al. found that the choice of standardization procedure can be more important for classification accuracy than the choice of the classifier itself [14].

Despite the demonstrated importance of distributional shift, methods to mitigate its impact are limited [13]. We address the problem directly, focusing on two sources of distributional shift in the context of sleep scoring: *nuisance variability*, caused by changes in the way signals are recorded, and *class balance variability*, caused by changes in the time spent in each sleep stage. Both forms of distributional shift might be present in a single dataset simultaneously. However, unlike nuisance variability, class balance variability should not be

removed because the primary motivation for sleep scoring is often to detect altered sleep behavior.

We use both a simple model and mouse *in vivo* data to demonstrate that standard  $z$ -scoring, which aims to remove nuisance variability, inappropriately reduces class balance variability. Therefore, classification algorithms using standard  $z$ -scoring as a preprocessing step will be biased towards underestimating changes in class balance relative to their training data, which limits their applicability in research settings where the fraction of time spent in each sleep stage is of interest. We developed a method, *mixture  $z$ -scoring*, for standardizing features of the EEG and EMG signals that disentangles nuisance and class balance variability using a small amount of labeled data for each subject. Because this requires brief but non-trivial user interaction with the recordings, we also present a free, open-source user interface that allows for rapid manual sleep scoring for standardization purposes followed by machine learning-based, automated sleep scoring. The software is available at <https://github.com/zekebarger/AccuSleep>.

## 2.3 Results

### 2.3.1 Mixture $z$ -scoring corrects for distributional shift in simulated data

In this section, we use a simulated experiment (Fig 2.2) to model distributional shift in its simplest form in order to illustrate the problems caused by standard  $z$ -scoring and resolved by mixture  $z$ -scoring. Empirical findings in the next section indicate that these results generalize to a more realistic setting.

The simulation is designed as follows: two experimental subjects, Subject I and Subject II, spend different amounts of time in each of two states, state 1 and state 2. In order to detect and quantify this difference, a data feature,  $\delta$ , is measured from each subject and paired observations of  $\delta$  and ground truth class labels from Subject I are used to train a logistic regression model. This model is then used to classify observations from Subject II.  $\delta$  is drawn from a mixture of Gaussians, with one Gaussian for each class. This scenario is analogous to an experiment where Subject I is a control mouse, Subject II is a mouse undergoing sleep deprivation, state 1 is wakefulness, state 2 is sleep, and  $\delta$  is the delta power in the EEG signal recorded from each mouse. The process of training the classifier on Subject I and then applying it on Subject II is analogous to the common practice of developing algorithms on wild-type populations but then applying them on wild-type and experimental subjects.

From Fig 2.2A, which shows kernel density estimates for the marginal and class-conditional distributions for both subjects, it can be seen that the distribution of  $\delta$  for Subject I differs from that for Subject II. When the distribution of the data on which a classifier is applied

differs from the training distribution, we say that there has been distributional shift. Class balance variability and nuisance variability are two major sources of distributional shift in scientific applications of classification algorithms.

Class balance variability is often the substance of the scientific inquiry that a classifier is meant to support. An experimental intervention, such as a drug, is expected to alter the amount of time spent in one or more sleep states and classification algorithms for sleep scoring are to be used to detect this change. The datasets in this simulated experiment have different class balances: Subject II spends 30% of the time in state 2 (blue), while Subject I spends 70% of the time in that state. This difference can be seen in the marginal distributions of  $\delta$  (black curves) for each subject.

In contrast, changes in the distribution of  $\delta$  due to nuisance variability should be removed so that observations from different subjects can be compared. Common sources of nuisance variability in the context of sleep scoring include the use of different types of recording equipment and different implantation sites for recording electrodes. The simplest form of nuisance variability is an affine transformation. In this example, the data from Subject II have the same class-conditional distributions as in Subject I, except for an affine transformation.

For an otherwise fixed distribution, an affine transformation can be undone by means of  $z$ -scoring. That is, the mean,  $\mu$ , and standard deviation,  $\sigma$ , of measurements  $\Phi$  are calculated for each subject and then values,  $Z$ , known as  $z$ -scores are computed:

$$Z = \frac{\Phi - \mu}{\sigma}. \quad (2.1)$$

If the distributions of measurements from two subjects differ only by an affine transformation, then after the application of  $z$ -scoring, they will be identical. Furthermore, they will have mean 0 and standard deviation 1, and so this procedure is also called *standardization*. It is typically beneficial for machine learning algorithms to operate on standardized features [15].

However, as panels B-F of Fig 2.2 demonstrate, when other sources of distributional shift are present, standard  $z$ -scoring is inappropriate. Standard  $z$ -scoring uses the marginal statistics of  $\delta$ , which are dependent on both nuisance and class balance variability. If we view the marginal distribution of  $\delta$  as a mixture of the class-conditional distributions for each class, its mean and variance are functions of the mixture weights (equivalently, the marginal probabilities of each class label) and the mean and variance of each class-conditional distribution. Its mean is given by the weighted average of the conditional means, while the variance is given by the law of total variance (Eq 2.13, Methods).

After standard  $z$ -scoring, the class-conditional distributions of  $\delta$  from the two subjects do not align as desired (compare the yellow curves in Fig 2.2B and C). The result is that a classifier trained to high performance on the  $z$ -scored data from Subject I will perform poorly on the  $z$ -scored data from Subject II (Fig 2.2D): a large fraction of observations are mislabeled as state 2 as the decision boundary aligns with the center of the distribution

for state 1. In this case, the introduced bias is opposite and almost equal to the effect of the change in class balance, leading to the misclassification of a significant fraction of the observations in Subject II as state 2 when they should be state 1. The resulting bias in the estimate of the effect size of the difference leads to a reduction in power and a false negative result for a bootstrapping test (Fig 2.2E). Because  $z$ -scoring is agnostic to class label, it is unable to disentangle the effect of the class balance from the effect of nuisance variability on the marginal statistics of  $\delta$ .

In order to remove nuisance variability but retain class balance variability, we introduce *mixture  $z$ -scoring*, an alternative form of  $z$ -scoring inspired by viewing the marginal distribution of the measurements as a mixture of class-conditional distributions. The mixture  $z$ -scored values  $Z_M$  corresponding to measured feature values  $\Phi$  are computed as:

$$Z_M = \frac{\Phi - w^\top \hat{\mu}}{\sqrt{w^\top (\hat{\sigma}^2 + (\hat{\mu} - w^\top \hat{\mu})^2)}} \quad (2.2)$$

where  $\hat{\mu}$  and  $\hat{\sigma}$  are vectors of label-conditioned means and standard deviations and  $w$  is a fixed vector that sums to 1. The value of the denominator and the value subtracted from  $\Phi$  in the numerator are the mixture  $z$ -scoring parameters, by analogy with the  $z$ -scoring parameters. Note that conditioning on labels means that mixture  $z$ -scoring requires some labeled data. If  $w$  is equal to the actual proportions of class labels for the values  $\Phi$ , then the mixture  $z$ -score parameters are the same as in standard  $z$ -scoring. When  $w$  is not equal to the class balance, then the mixture  $z$ -score parameters are equal to what the  $z$ -score parameters would have been, had the class balance of  $\Phi$  been  $w$ . Mixture  $z$ -scoring thus removes any effect of affine nuisance variability on the marginal distribution of  $\delta$  while preserving any effect of class balance variability. We refer to this process as *mixture standardization*, by analogy with  $z$ -score standardization. See the Methods section for further details.

The results of using mixture  $z$ -scoring are shown in Fig 2.2, panels F-I. Panels F and G show the marginal and class-conditional distributions of mixture-standardized  $\delta$  for Subjects I and II, respectively. Note that the class-conditional distributions are aligned, unlike in panels B and C. The result is that a classifier trained to high accuracy on the mixture-standardized data from Subject I also performs well on the data from Subject II when it is mixture-standardized with the same weights (panel H). A bootstrapping test comparing the time spent in state 2 for the two subjects now returns a true positive (panel I).

This example demonstrates that standard  $z$ -scoring of features used as inputs to a classifier can, in a simple case, lead to systematic classification errors when the class balance of the data on which the classifier is applied differs from the balance during training. This systematic error can lead to incorrect scientific conclusions. However, mixture  $z$ -scoring, which standardizes data without removing class balance variability, substantially reduces the error rate.



### 2.3.2 Mixture $z$ -scoring reduces bias when classifying mouse *in vivo* data

While the above section indicates that mixture  $z$ -scoring can correct for distributional shift in a simulated experiment with a simple classifier, it remains to be seen whether distributional shift still poses a problem in a more realistic setting with state-of-the-art methods. To address this question, we applied two automated sleep scoring algorithms to mouse EEG and EMG recordings:

#### 1. SPINDLE

This algorithm, from [10], comprises a convolutional neural network (CNN) and a hidden Markov model (HMM). The CNN operates on multi-channel EEG and EMG spectrograms and comprises max-pooling, convolution, and max-pooling followed by two fully-connected layers. It has 6.8M parameters, distributed primarily in the first fully-connected layer. The HMM is used to constrain the transitions between sleep stages predicted by the CNN. See [10] for details. EMG activity and individual frequency bands of the spectrogram are log-transformed and, importantly,  $z$ -scored on a per-recording basis.

#### 2. Sleep Scoring Artificial Neural Network (SS-ANN)

We implemented a simple CNN with fewer than 20K learnable parameters that, similar to SPINDLE, operates on log-transformed EEG spectrograms and EMG activity. We refer to this network as SS-ANN. It uses three convolution-ReLU-maxpool modules with batch normalization, followed by a linear classifier. See the Methods section for details of the network architecture and a list of the datasets used for training and testing SS-ANN in each experiment.

Each algorithm was applied to recordings from three wild-type mice. We applied SPINDLE to recordings from Cohort A in [10]. We applied SS-ANN to a separate set of recordings which we collected and scored (see Methods). This ensured that any observed bias in either method would not be a result of differences between the experts who scored the training datasets.

To introduce class balance variability and produce a controlled simulation of the effect of experimental manipulations, we programmatically varied the amounts of NREM sleep, REM sleep, and wakefulness in each recording (Fig 2.3A). Rebalancing was achieved by randomly removing bouts of each sleep stage until a specified balance was reached. To preserve the natural temporal structure of the bouts as much as possible, we included at least eight seconds of NREM before, and four seconds of wakefulness after, each REM bout. We considered class balances in the range of 5-95% wakefulness, 5-95% NREM, and 0-25% REM because class balances outside this range are unlikely to be observed in practice.

In order to demonstrate that any observed biases are general to  $z$ -scoring, rather than specific to SPINDLE, and that mixture  $z$ -scoring can eliminate them, we trained two different versions of SS-ANN. For one, the inputs during both training and testing were preprocessed with standard  $z$ -scoring. For the other, both inputs were preprocessed using mixture  $z$ -scoring with weight vector  $w$  given by the class balance in the training set. Mixture  $z$ -scoring requires a small amount of labeled data from each class for each subject, which we took from separate recordings from the three mice in this dataset.

The results of these numerical experiments are shown in Fig 2.3. When using standard  $z$ -scoring, both algorithms showed a classification bias that typically made the aggregate label distribution look more like the class balance of their training datasets (Fig 2.3B,C). The magnitude of the classification error increased as the class balance was shifted further from the balance of the training data (Fig 2.3E). These results indicate that this bias occurs across classifiers. However, when mixture  $z$ -scoring was used as a preprocessing step, there was a dramatic reduction in bias: estimated label fractions were close to the true fractions in all cases (Fig 2.3E), eliminating the contraction towards the class balance of the training data (Fig 2.3D).

We quantified the bias in estimation by measuring the total variation distance between the ground truth label fractions and the label fractions estimated by each classifier (Fig 2.3F). The total variation distance,  $\delta$ , between two distributions  $q$  and  $q'$  is  $\delta(q, q') = \|q - q'\|_1$ . We found that there was no significant difference between the two algorithms when standard  $z$ -scoring was used (mean for SPINDLE: 0.24, mean for SS-ANN: 0.25, Student's  $t=-0.22$ ,  $p = 0.83$ ). We further found that both SPINDLE and SS-ANN had substantially greater bias when using standard  $z$ -scoring than did SS-ANN using mixture  $z$ -scoring (mean for SS-ANN with mixture: 0.04, SPINDLE vs SS-ANN mixture  $t=8.19$ ,  $p \ll 0.01$ ; SS-ANN standard vs SS-ANN mixture  $t=10.33$ ,  $p \ll 0.01$ ). These results demonstrate that while standard  $z$ -scoring introduces substantial bias in state-of-the-art machine learning classifiers in cases when class balance variability is present, methods such as mixture  $z$ -scoring can significantly reduce this bias.

### 2.3.3 Validation of SS-ANN

To demonstrate the utility of SS-ANN paired with mixture  $z$ -scoring for automated sleep scoring, we evaluated their performance using several different metrics (Fig 2.4). On held-out test data, we found that SS-ANN achieved 96.8% accuracy against expert annotations, comparable to the range of values reported for inter-scorer agreement (typically in the range of 90-96%) [10, 16, 17] and to the performance of SPINDLE (95-96%, see Methods), despite SS-ANN having over 300 times fewer parameters. Agreement was also high for each class individually, as evidenced by the high values on the diagonal of the confusion matrix and low values off of it (Fig 2.4A). The receiver operating characteristic (ROC) curves for each class were far away from the unity line (bootstrap  $p$ -value  $\ll 0.001$ ) (Fig 2.4B).

The experiment in Fig 2.3 demonstrated that SS-ANN can generalize across simulated experimental conditions that create different class balances. To determine whether it can also generalize across subjects, we trained and tested the network on data from each individual subject in our cohort (Fig 2.4C). Classification accuracy was comparable to rates of inter-scoring agreement for all train-test pairs despite the reduced size of the training sets, indicating good generalization.

Finally, we investigated the relationship between classification accuracy and the amount of labeled data used for mixture  $z$ -scoring. As described in the Methods section, mixture  $z$ -scoring requires labeled data from each subject in order to estimate a mean and variance for each data feature within each class. Though the class balance of the labeled sample has no direct effect on the estimation of these parameters, it is important to determine how many labeled epochs are required to attain accurate enough parameter estimates to support classification.

We held out a number of labeled epochs chosen at random from EEG/EMG recordings and used these to perform mixture  $z$ -scoring followed by automatic classification of the remaining epochs. Classification accuracy increased with larger held-out portions until reaching a plateau at approximately 10 minutes of labeled data (Fig 2.4D). Note that performance is already high for even small sample sizes, under one minute. This should not be taken to indicate that mixture  $z$ -scoring is unnecessary to obtain high performance, since the animals in our cohort did not undergo any experimental manipulations and the intent of mixture  $z$ -scoring is to account for the possibility of observing altered class balances. Additionally, note that small sample sizes result in uncorrelated errors in the estimation of the  $z$ -score parameters for each feature, which have less impact on classification error than do the correlated errors caused by class balance variability.

### 2.3.4 AccuSleep: free, open-source software for automated sleep scoring

The benefit of mixture  $z$ -scoring for generalization across class balances comes with the requirement that for each subject, there must be enough labeled epochs of each class to estimate the class-conditional means and variances of each feature. In practice, labeling these epochs requires a user to interact with the data in a non-trivial way. We addressed this issue by creating AccuSleep, a set of MATLAB graphical user interfaces that allow for manual scoring of EEG/EMG data (Fig 2.5) followed by automated scoring using SS-ANN (Fig 2.6). The mixture  $z$ -scoring parameters, once calculated for a given subject, can be used to score other recordings from the same subject. The workflow for scoring recordings is simple:

1. **Select EEG and EMG data.** In the interface shown in Fig 2.6, the user selects the data files, enters the sampling rate and epoch length, and sets the output location for

all recordings from a given subject. Arbitrary epoch lengths can be used, provided that the neural network used for automatic scoring was trained on data scored at the same temporal resolution.

2. **Set mixture  $z$ -scoring parameters.** If mixture  $z$ -scoring parameters have already been calculated for this subject, they can be loaded. If not, the user scores a small number of epochs of each state manually using the interface shown in Fig 2.5. AccuSleep then calculates and saves the parameters.
3. **Load a trained copy of SS-ANN.** The user loads a network for automated classification. The trained network validated in Fig 2.4 is included with the software, along with MATLAB functions to retrain the network on new data.
4. **Score sleep automatically.** After validating the inputs, AccuSleep uses SS-ANN to perform automatic sleep scoring.

For an experienced user, manually scoring the number of epochs required for maximum classification accuracy (Fig 2.4D) requires roughly 2 minutes, far shorter than it would take to score all of the data from that subject. Thereafter, scoring of additional data from the same subject does not require any manual labeling, preserving the scaling and efficiency benefits of automated scoring.

## 2.4 Discussion

Supervised machine learning is well suited to the task of sleep scoring: labeled data are plentiful, and contemporary algorithms can learn from minimally processed EEG and EMG data to achieve classification accuracy comparable to inter-rater reliability. Nevertheless, machine learning algorithms are still not widely used for sleep scoring in research. We suspect there are two reasons for this: low usability, since applying machine learning methods can require specialized knowledge or skills; and poor generalization, since variability in the EEG and EMG signals due to inter-subject and inter-laboratory differences, or distributional shift, poses a challenge to the generalization of any algorithm. Both of these issues must be addressed if machine learning is to be widely adopted for sleep scoring.

We propose mixture  $z$ -scoring (Eq 2.2) as a solution to the problem of generalization posed by distributional shift due to simultaneous nuisance and class balance variability. Standard  $z$ -scoring serves this purpose well when class balance variability is low, but not when class balance variability is high, as it often is in scientific settings. In our experiments, which simulated changes in class balance such as might occur in a sleep study, we found that a classifier using standard  $z$ -scoring as a preprocessing step performed poorly on data from a subject with a different class balance than its training set (Fig 2.3B-E). This can produce undesirable results in a research setting because the effect of a manipulation that changes

the quantity of sleep or wakefulness would be poorly estimated (Fig 2.2F). Mixture  $z$ -scoring captures the two sources of variability independently, removing only the former. On both simulated (Fig 2.2E,I) and real (Fig 2.3B-F) data with artificially-varied class balance, our method improves generalization and estimation of class balance. We expect that the same is true for real data with natural class balance variability—for example, recordings collected during and after a sleep deprivation protocol where the fractions of wakefulness and sleep are both different from baseline conditions.

We also introduce a new classification algorithm for rodent sleep scoring, SS-ANN (Sleep Scoring Artificial Neural Network). This convolutional neural network achieves comparable accuracy to inter-scorer agreement and to another neural network method, SPINDLE [10] (SS-ANN: 96.8%; SPINDLE: 94.8-96.2%, see Methods). This is achieved with 300x fewer parameters, at a higher temporal resolution, and without using a Hidden Markov Model to constrain transitions between sleep states. The vast majority of parameters in both networks are in the fully-connected layer that follows the convolutional component of the architecture. Thus, the difference in parameter count comes from SS-ANN’s greater use of convolution and pooling. The possibility that even simpler models with fewer parameters might achieve comparable performance remains to be explored.

While mixture  $z$ -scoring dramatically reduces bias due to class balance variability and improves generalization, this comes at a cost: a sample of labeled data from each subject must be provided in order to capture subject-specific variability. Unsupervised, or label-free, methods for handling class balance variability would avoid this requirement, but have other costs. In [17], certain feature values are explicitly assumed to be an ordered mixture-of-Gaussians, allowing for the threshold of a linear classifier to be placed without requiring labels. Mixture  $z$ -scoring does not make any assumption about the shape of the class-conditional feature distributions, nor does it rely on the specific form of the downstream classifier. The two methods are equivalent in the case of mixture-of-Gaussians data and a linear classifier. Alternative unsupervised methods that do not involve distributional assumptions might include data augmentation with inputs whose nuisance and class balance variability are altered programmatically or randomly, so long as these augmentation methods are applied upstream of the preprocessing step. One advantage of our supervised method over this alternative is that it applies to any class balance and any affine nuisance variability, rather than only for ranges included in the augmentation step. Further, our method preserves any affine variability that is useful to classification, while data augmentation removes it. For example, applying data augmentation to the linear classification problem in Fig 2.2 would result in a classifier with poor performance. The success of threshold-based methods [3, 17] indicates that affine features, such as EMG power, are very useful for classification, suggesting that data augmentation would be harmful to performance. Finally, our results indicate that the amount of labeled data required to achieve high accuracy is on the order of minutes (Fig 2.4D), indicating that the time cost of the labeling step is quite small.

To preserve the speed and scalability benefits of automated scoring that includes mix-

ture  $z$ -scoring, and therefore some manual labeling, as a preprocessing step, we aimed to make the labeling step as streamlined as possible. To this end, we developed AccuSleep: a free, open-source MATLAB package that provides graphical interfaces for manual and deep learning-based, automated sleep scoring (Fig 2.6). Within AccuSleep, polysomnographic recordings can be manually scored to provide the labeled data for mixture  $z$ -scoring (Fig 2.5). These mixture parameters can then be used to score all recordings from the same subject automatically. Automatic classification is performed using a copy of SS-ANN trained on the data collected for these experiments (see Methods).

While SS-ANN showed good generalization across the 10 mice in our cohort (Fig 2.4C), it is possible that recordings collected from mice with different genotypes or differently placed electrodes would have substantially different EEG spectra from those in our dataset. To account for this possibility, we leverage one of the advantages of end-to-end learning—the fact that training new models is simple given labeled data—by including a module in AccuSleep that can train a new version of SS-ANN based on a sample of labeled data. Such a model could be provided alongside a research study to increase the replicability of its sleep scoring methodology. The time and computational resources required for this training process are minimal, owing to the small number of parameters in SS-ANN.

In summary, we developed a standardization procedure, mixture  $z$ -scoring, that simultaneously corrects for distributional shift due to both nuisance variability and changes in class balance. We demonstrated that our method improves the generalization of sleep scoring algorithms and provide software to enable its application in sleep research. We expect that this software, available at <https://github.com/zekebarger/AccuSleep>, will be useful to the research community. More broadly, we note that data standardization that goes beyond textbook  $z$ -scoring and accounts for class balance changes across experimental units is ubiquitous in experimental sciences (e.g.,  $z$ -scoring of electrocorticography recordings from humans [18]). Such methods are less commonly used in a machine learning setting, where algorithms are typically formulated assuming no distributional shift and validated using test sets that have similar class balance to the training set, such as held-out data. As demonstrated here, the practice of training algorithms in this manner on one type of particularly convenient experimental subject and then applying them on another, ignoring class balance variation, can lead to incorrect scientific conclusions. This suggests that mixture  $z$ -scoring could improve the accuracy of machine learning algorithms across scientific domains.

## 2.5 Methods

### 2.5.1 Polysomnographic recordings

All experimental procedures were approved by the Animal Care and Use Committee at the University of California, Berkeley. Animals were housed on a 12-hour dark/12-hour light cycle (light on between 7:00 and 19:00). Adult C57BL/6 mice (10-20 weeks old) were anes-

thetized with 1.5%–2% isoflurane and placed in a stereotaxic frame. Body temperature was kept stable throughout the procedure using a heating pad. After asepsis, the skin was incised to expose the skull, and the overlying connective tissue was removed. For EEG and EMG recordings, a reference screw was inserted into the skull on top of the right cerebellum. EEG recordings were made from two screws on top of the left and right cortex, at anteroposterior  $-3.5$  mm and medio-lateral  $\pm 3$  mm. Two EMG electrodes were inserted into the neck musculature. Insulated leads from the EEG and EMG electrodes were soldered to a pin header, which was secured to the skull using dental cement. All efforts were made to minimize suffering during and after surgery.

Recordings were made with the mice in their home cages placed in sound-attenuating boxes. Five 4-hour recordings were collected from each of 10 mice, and two 24-hour recordings were collected from five of those mice. For 24-hour recordings, recording started at 19:00 following 24 hours of habituation and lasted 48 hours. For four-hour recordings, recording started at 13:00 following two hours of habituation. The pin header was connected to a flexible recording cable via a mini-connector. Signals were recorded with a TDT RZ5 amplifier for the 24-hour recordings (bandpass filter, 1-750 Hz; sampling rate, 1,500 Hz) or an Intan Technologies RHD-2132 amplifier for the 4-hour recordings (bandpass filter, 1-500 Hz; sampling rate, 1,000 Hz). Sleep stages were scored manually in 2.5-second epochs by an expert scorer according to standard criteria. The complete dataset is available at <https://osf.io/py5eb/>.

The data used for training and testing SS-ANN in each experiment are described below. Unless otherwise specified, SS-ANN was trained using three 4-hour recordings from each of the 10 mice.

- Fig 2.3: tested on three programmatically rebalanced 12-hour light cycle recordings (extracted from the 24-hour recordings) from three mice
- Fig 2.4A,B: tested on two held-out 4-hour recordings from each of 10 mice
- Fig 2.4C, off-diagonal: trained on all five 4-hour recordings from each mouse, tested on all five 4-hour recordings from each of the other nine mice
- Fig 2.4C, on-diagonal: five-fold cross-validation using all five 4-hour recordings from each mouse
- Fig 2.4D: as for Fig 2.4A,B, but with a held-out portion of each test recording used for mixture  $z$ -scoring

Three mice were not used for recordings due to low signal-to-noise ratio (SNR) in the EEG and EMG signals. One mouse was excluded from our analyses because the SNR of its signals decreased before data collection was completed.

## 2.5.2 Sleep scoring algorithm

### Data preprocessing

EEG and EMG signals were downsampled to 128 Hz. We used the Chronux toolbox [19] to calculate a multi-taper spectrogram of the EEG signal between 0-50 Hz with a 5 second window and 2.5 second step. We downsampled by a factor of 2 between 20-50 Hz to reduce the number of parameters in the classifier. To calculate EMG activity, we bandpass filtered the EMG signal between 20-50 Hz and took the root-mean-square of the signal in each epoch. To build the complete feature set for each recording, we concatenated the EEG spectrogram with 9 copies of the EMG activity. Since the spectrogram has 176 frequency components, each recording becomes a  $185 \times n$  matrix with 185 features for each of  $n$  epochs.

### SS-ANN architecture

The inputs to SS-ANN were  $185 \times 13$  pixel grayscale images, representing 32.5-second periods of the standardized joint EEG/EMG spectrogram centered on each epoch. We created a basic CNN architecture using the MATLAB Statistics and Machine Learning Toolbox (MATLAB, The MathWorks, Natick, MA): 3 convolution - batch normalization - ReLU - max pooling modules, followed by a fully connected layer, softmax layer, and classification layer. The convolution layers had filter size 3 with 8, 16, and 32 filters per layer. The max pooling layer had size 2 and stride 2. The network was trained using stochastic gradient descent with momentum and a mini-batch size of 256 for 10 epochs. The learning rate was 0.015, reduced by 15% each epoch. Classes were balanced prior to training by randomly oversampling the classes with the fewest examples to reach the number of examples in the largest class. Following the classification step, sleep stages were refined by assigning bouts shorter than 5 seconds to the surrounding stage.

### Comparison between SS-ANN and SPINDLE

Up to 20% of epochs in the recordings used for training or testing SPINDLE were scored as artifacts [10], but inter-rater agreement for artifact detection was low (in the range of 20-30%). The accuracy of SPINDLE reported in [10] was calculated only using epochs not labeled as artifacts and on which two expert raters agreed. These criteria may remove some of the most difficult-to-classify epochs. We used the SPINDLE online service with artifact detection disabled to re-score the datasets used in that publication, obtaining accuracies of 96.2%, 95.1%, 94.8% on Cohorts A, B, and C versus the labels of Expert 1.

### Mixture $z$ -scoring

For a given set of  $n + N$  observations of feature values  $X$ ,  $n$  of which have paired observations of labels  $L$ , and fixed choice of baseline mixture weights  $w$ , we perform *mixture  $z$ -scoring* to



obtain standardized observations  $Z_M$  as follows:

$$Z_M = \frac{X - w^\top \hat{\mu}}{\sqrt{w^\top (\hat{\sigma}^2 + s)}}$$

defining subtraction and division between vectors and scalars and squaring of vectors as the element-wise versions of their scalar equivalents, where the vectors  $\hat{\mu}$ ,  $\hat{\sigma}$ , and  $s$  are as below, denoting by  $X_l$  the set of observations with paired label equal to  $l$ :

$$\hat{\mu}_l = \frac{1}{n} \sum_{x \in X_l} x \quad (2.3)$$

$$\hat{\sigma}_l^2 = \frac{1}{n} \sum_{x \in X_l} (x - \hat{\mu}_l)^2 \quad (2.4)$$

$$s_l = (\hat{\mu} - w^\top \hat{\mu}) \quad (2.5)$$

We found that choosing the baseline mixture weights  $w$  to be close to the class prevalences in a reference dataset worked well.

This algorithm can be motivated as follows: let  $\phi$  be a random variable corresponding to the feature value and  $Y$  be the random variable corresponding to the class label. We can break down the distribution of the feature value,  $P(\phi)$ , into a mixture of the distributions conditioned on the class label,  $P(\phi|Y = i)$ , each with a corresponding mean  $\mu_i$  and variance  $\sigma_i^2$

$$w_i := P(Y = i) \quad (2.6)$$

$$\mathbb{E}[\phi|Y = i] := \mu_i \quad (2.7)$$

$$\sigma_i^2 := \mathbb{V}[\phi|Y = i] \quad (2.8)$$

The overall mean  $\mu_G$  and variance  $\sigma_G^2$  of this mixture distribution can be written in terms of the means and variances of its components as follows, writing  $\mu$  and  $\sigma^2$  for the vectors of class-conditional means and variances, and  $w$  for the vector of class probabilities:

$$\mu_G := \mathbb{E}[\phi] = w^\top \mu \quad (2.9)$$

$$\sigma_G^2 := \mathbb{V}[\phi] = w^\top (\sigma^2 + s) \quad (2.10)$$

$$\mathbb{V}[\phi] = w^\top \sigma^2 + w^\top (\mu - w^\top \mu)^2 \quad (2.11)$$

$$s := (\mu - w^\top \mu)^2 \quad (2.12)$$

where the expression for  $\mu_G$  comes from the linearity of the expectation and the expression for  $\sigma_G^2$  comes from the law of total variance:

$$\mathbb{V}[\phi] = \mathbb{E}[\mathbb{V}[\phi|Y]] + \mathbb{V}[\mathbb{E}[\phi|Y]] \quad (2.13)$$

The vector  $s$  is analogous to a “between sum-of-squares” in an ANOVA.

In this representation, we can write the probabilistically  $z$ -scored version of  $\phi$ , which we write  $\phi_Z$ , as:

$$\phi_Z = \frac{\phi - \mu_G}{\sigma_G} = \frac{\phi - w^\top \mu}{\sqrt{w^\top (\sigma^2 + s)}} \quad (2.14)$$

The key utility of this representation is that it separates out contributions to the overall mean and variance by the means and variances of each group from contributions to the overall mean and variance from the *prevalence* of each group. Note that  $w$ ,  $\mu$ ,  $\sigma^2$ , and  $s$  all typically need to be estimated from data, which results in the typical, exact form of  $z$ -scoring.

Now instead suppose we observe a nuisance affected version,  $\tilde{\phi}$ , with the same class balance  $w$ . We presume the nuisance variability acts to, in expectation, scale and shift the distribution of  $\phi$  by scaling parameter  $a$  and shift parameter  $b$ :

$$\tilde{\mu} := \mathbb{E} [\tilde{\phi} | Y = i] = a \mathbb{E} [\phi | Y = i] + b \quad (2.15)$$

$$\tilde{\sigma}^2 \phi \mathbb{V} [\tilde{\phi} | Y = i] = a^2 \mathbb{V} [\phi | Y = i] \quad (2.16)$$

$$\tilde{s} := (\tilde{\mu} - w^\top \tilde{\mu})^2 \quad (2.17)$$

The method for probabilistic  $z$ -scoring remains the same, with the new group means and variances substituted in:

$$\tilde{\phi}_Z = \frac{\tilde{\phi} - \tilde{\mu}_G}{\tilde{\sigma}_G} = \frac{\tilde{\phi} - w^\top \tilde{\mu}}{\sqrt{w^\top (\tilde{\sigma}^2 + \tilde{s})}} \quad (2.18)$$

That is, in the absence of changes in class balance, we can remove affine nuisance variability by subtracting off a weighted sum of the class-conditional means and dividing by a weighted sum of the class-conditional variances and the weighted variability of the means.

This suggests a method of  $z$ -scoring to remove affine nuisance variability for observations,  $\hat{\phi}$ , with different class balances  $\hat{w}$ . We compute the equivalent  $\hat{\mu}$  and  $\hat{\sigma}^2$ , which are class-specific statistics, then plug them into Eq 2.14 with the weights given by  $w$  instead of  $\hat{w}$ :

$$\hat{\phi}_{Z_m} = \frac{\hat{\phi} - w^\top \hat{\mu}}{\sqrt{w^\top (\hat{\sigma}^2 + (\hat{\mu} - w^\top \hat{\mu})^2)}} \quad (2.19)$$

The result, applied to a finite dataset, is the operation defined by Eq 2.2. In essence: we try to perform the same nuisance variability-removing  $z$ -scoring operation we would have done, had there not been a class balance change.

Note that in general the resulting data  $Z_m$  no longer has mean 0 or standard deviation 1, since that is only true when  $\hat{w} = w$ .

## Acknowledgements

We would like to thank Chak Foon Tso for his assistance testing the software and Hayley Bounds for her assistance designing the neural network.

## References

- [1] Stephenson, R., Caron, A. M., Cassel, D. B., & Kostela, J. C. (2009). Automated analysis of sleep–wake state in rats. *Journal of Neuroscience Methods*, *184*(2), 263–274. <https://doi.org/10.1016/j.jneumeth.2009.08.014>
- [2] Kohtoh, S., Taguchi, Y., Matsumoto, N., Wada, M., Huang, Z., & Urade, Y. (2008). Algorithm for sleep scoring in experimental animals based on fast fourier transform power spectrum analysis of the electroencephalogram. *Sleep and Biological Rhythms*, *6*(3), 163–171. <https://doi.org/10.1111/j.1479-8425.2008.00355.x>
- [3] Bastianini, S., Berteotti, C., Gabrielli, A., Vecchio, F. D., Amici, R., Alexandre, C., Scammell, T. E., Gazea, M., Kimura, M., Martire, V. L., Silvani, A., & Zoccoli, G. (2014). SCOPRISM: A new algorithm for automatic sleep scoring in mice. *Journal of Neuroscience Methods*, *235*, 277–284. <https://doi.org/10.1016/j.jneumeth.2014.07.018>
- [4] Kreuzer, M., Polta, S., Gapp, J., Schuler, C., Kochs, E., & Fenzl, T. (2015). Sleep scoring made easy—semi-automated sleep analysis software and manual rescoring tools for basic sleep research in mice. *MethodsX*, *2*, 232–240. <https://doi.org/10.1016/j.mex.2015.04.005>
- [5] Gross, B. A., Walsh, C. M., Turakhia, A. A., Booth, V., Mashour, G. A., & Poe, G. R. (2009). Open-source logic-based automated sleep scoring software using electrophysiological recordings in rats. *Journal of Neuroscience Methods*, *184*(1), 10–18. <https://doi.org/10.1016/j.jneumeth.2009.07.009>
- [6] Shantilal, Donohue, K. D., & O’Hara, B. F. (2008). SVM for automatic rodent sleep–wake classification, In *IEEE SoutheastCon 2008*.
- [7] Crisler, S., Morrissey, M. J., Anch, A. M., & Barnett, D. W. (2008). Sleep-stage scoring in the rat using a support vector machine. *Journal of Neuroscience Methods*, *168*(2), 524–534. <https://doi.org/10.1016/j.jneumeth.2007.10.027>
- [8] Rempe, M. J., Clegern, W. C., & Wisor, J. P. (2015). An automated sleep-state classification algorithm for quantifying sleep timing and sleep-dependent dynamics of electroencephalographic and cerebral metabolic parameters. *Nature and Science of Sleep*, *7*, 85–99. <https://doi.org/10.2147/NSS.S84548>
- [9] Vilamala, A., Madsen, K. H., & Hansen, L. K. (2017). Deep convolutional neural networks for interpretable analysis of EEG sleep stage scoring. *arXiv e-prints*, arXiv 1710.00633.
- [10] Miladinović, D., Muheim, C., Bauer, S., Spinnler, A., Noain, D., Bandarabadi, M., Gallusser, B., Krummenacher, G., Baumann, C., Adamantidis, A., Brown, S. A., & Buhmann, J. M. (2019). SPINDLE: End-to-end learning from EEG/EMG to extrap-

- olate animal sleep scoring across experimental settings, labs and species. *PLOS Computational Biology*, 15(4), 1–30. <https://doi.org/10.1371/journal.pcbi.1006968>
- [11] Supratak, A., Dong, H., Wu, C., & Guo, Y. (2017). DeepSleepNet: a model for automatic sleep stage scoring based on raw single-channel EEG. *arXiv e-prints*, arXiv 1703.04046, arXiv:1703.04046.
- [12] Schwabedal, J. T. C., Sippel, D., Brandt, M. D., & Bialonski, S. (2018). Automated classification of sleep stages and EEG artifacts in mice with deep learning. *arXiv e-prints*, arXiv 1809.08443, arXiv:1809.08443.
- [13] Quiñonero-Candela, J., Sugiyama, M., Schwaighofer, A., & Lawrence, N. D. (Eds.). (2008). *Dataset Shift in Machine Learning (Neural Information Processing series)*. The MIT Press.
- [14] Katsageorgiou, V., Lassi, G., Tucci, V., Murino, V., & Sona, D. (2015). Sleep-stage scoring in mice: The influence of data pre-processing on a system’s performance, In *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. <https://doi.org/10.1109/EMBC.2015.7318433>
- [15] LeCun, Y. A., Bottou, L., Orr, G. B., & Müller, K.-R. (2012). Efficient backprop. In G. Montavon, G. B. Orr, & K.-R. Müller (Eds.), *Neural networks: Tricks of the trade: Second edition* (pp. 9–48). Springer Berlin Heidelberg. [https://doi.org/10.1007/978-3-642-35289-8\\_3](https://doi.org/10.1007/978-3-642-35289-8_3)
- [16] Yaghouby, F., O’Hara, B. F., & Sunderam, S. (2016). Unsupervised estimation of mouse sleep scores and dynamics using a graphical model of electrophysiological measurements. *International Journal of Neural Systems*, 26(04), 1650017. <https://doi.org/10.1142/S0129065716500179>
- [17] Bagur, S., Lacroix, M. M., de Lavilléon, G., Lefort, J. M., Geoffroy, H., & Benchenane, K. (2018). Harnessing olfactory bulb oscillations to perform fully brain-based sleep-scoring and real-time monitoring of anaesthesia depth. *PLOS Biology*, 16(11), 1–32. <https://doi.org/10.1371/journal.pbio.2005458>
- [18] Bouchard, K. E., Mesgarani, N., Johnson, K., & Chang, E. F. (2013). Functional organization of human sensorimotor cortex for speech articulation. *Nature*, 495(7441), 327–332. <https://doi.org/10.1038/nature11911>
- [19] Bokil, H., Andrews, P., Kulkarni, J. E., Mehta, S., & Mitra, P. P. (2010). Chronux: A platform for analyzing neural signals. *Journal of Neuroscience Methods*, 192(1), 146–151. <https://doi.org/10.1016/J.JNEUMETH.2010.06.020>

Figure 2.1: **Overview of the signal collection process for sleep scoring in mice.** A: schematic of EEG and EMG recordings. An EEG electrode is inserted over the hippocampus, a reference electrode is placed in the cerebellum, and an EMG electrode is inserted into the neck musculature. B: sample EEG and EMG recordings. Scale bar: 1 s, 0.25 mV. C: example EEG spectrograms and root-mean-square EMG activity for each sleep stage.

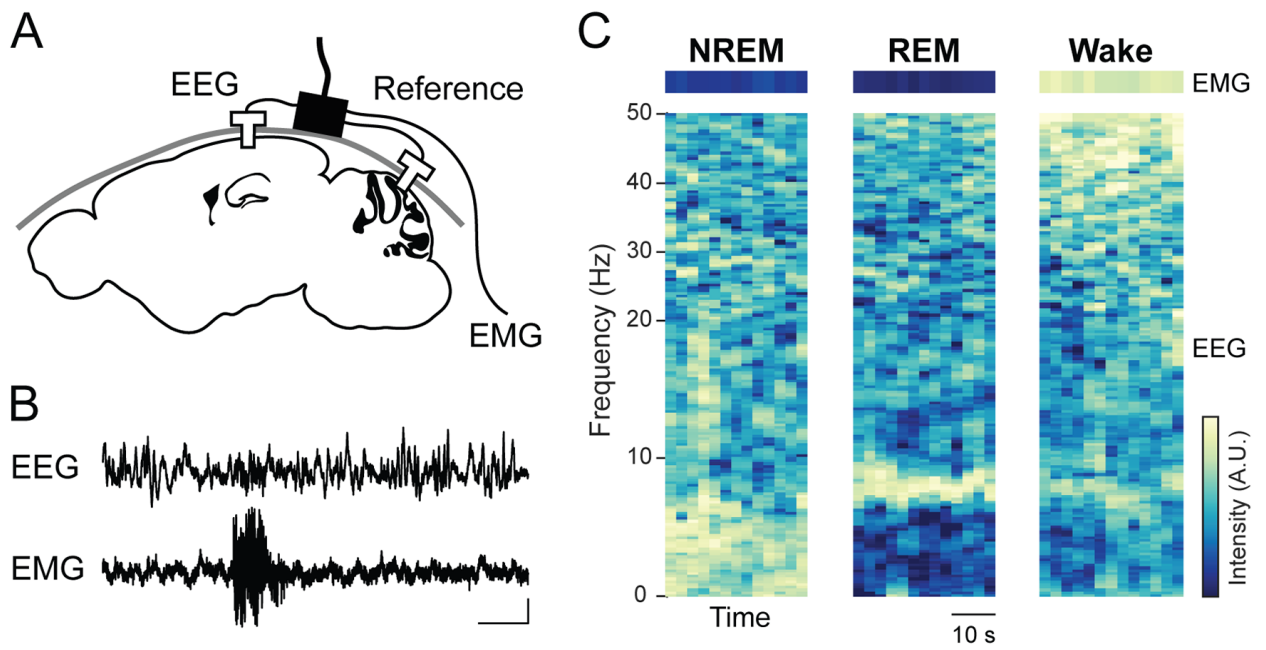


Figure 2.1. Overview of the signal collection process for sleep scoring in mice.

Figure 2.2: **Correcting for distributional shift prevents a false negative in a simple model** A: marginal and class-conditional distributions for a feature,  $\delta$ , recorded from synthetic Subject I (dashed lines), which is in state 2 70% of the time, and from synthetic Subject II (solid lines), which is in state 2 30% of the time. Marginal distributions are in black, the distribution for state 2 in yellow, state 1 in blue. The distributions differ in class balance and by an affine shift. Left column, B-E: the results of applying standard  $z$ -scoring to this synthetic data. The  $x$ -axis, representing the value of  $\delta$ , is shared across B-D. In B and C, the marginal and class-conditional distributions are plotted after standard  $z$ -scoring. D: the output of a logistic classifier trained on the  $z$ -scored data from Subject I (in color; decision threshold represented by gray bar) compared to the ground-truth state (on the  $y$ -axis) for a selection of data points from Subject II. The  $y$ -values are jittered to improve legibility. E: the estimated fraction of time in state 2 using labels given by the classifier from D. Error bars show approximate 95% confidence intervals (CIs) for this fraction, obtained by boot-strapping. The symbol “n.s.” indicates that the estimated fraction for Subject II fell within the 95% CI for Subject I. Right column, F-I: as B-E, but using mixture  $z$ -scoring to correct for distributional shift. Each panel on the right-hand-side shares its  $y$ -axis with the matching panel on the left-hand side. The \* symbol indicates that the estimated fraction for Subject II fell outside the 95% CI for Subject I.

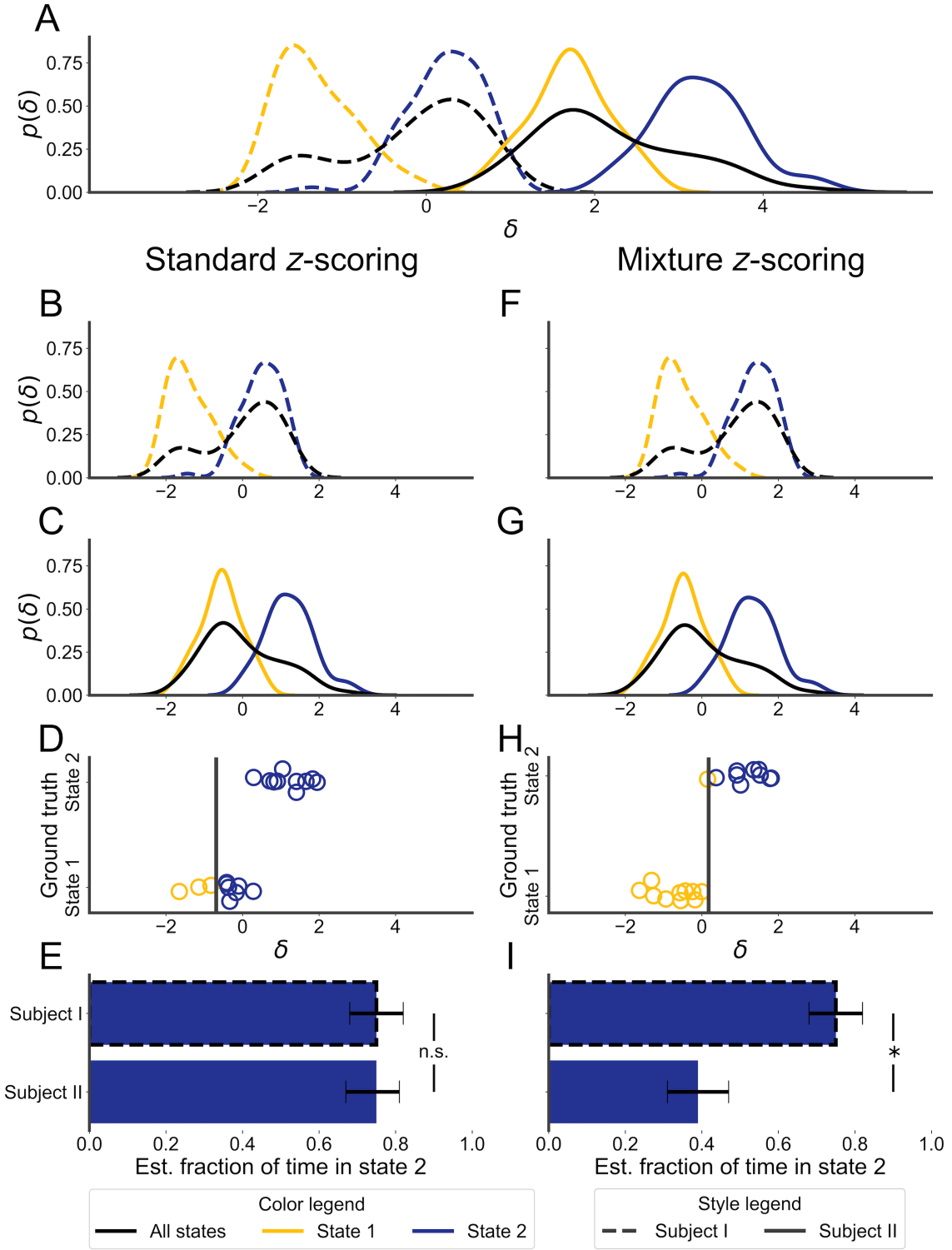


Figure 2.2. Correcting for distributional shift prevents a false negative in a simple model.



Figure 2.3: **Comparison of sleep scoring algorithms on recordings with programmatically varied class balances.** EEG and EMG recordings were programmatically altered to have different proportions of each sleep stage. A-D: Each marker indicates the amount of wakefulness and NREM in recordings with a given class balance, averaged across three mice (one recording per mouse). A: class balances of the recordings according to ground truth manual labels. B: predictions by SS-ANN with standard  $z$ -scoring for each recording. Star indicates the class balance in the training dataset for SS-ANN. C: as B, for SPINDLE. Star indicates class balance of training data for SPINDLE. D: as B, for SS-ANN with mixture  $z$ -scoring. E: predicted amounts of wakefulness in recordings containing 5% REM, according to each algorithm. F: mean total variation distance between class balances of algorithmic predictions and manual labels. Error bars show SEM. Asterisk indicates  $p$ -value  $< 0.001$  (Student's  $t$ -test).

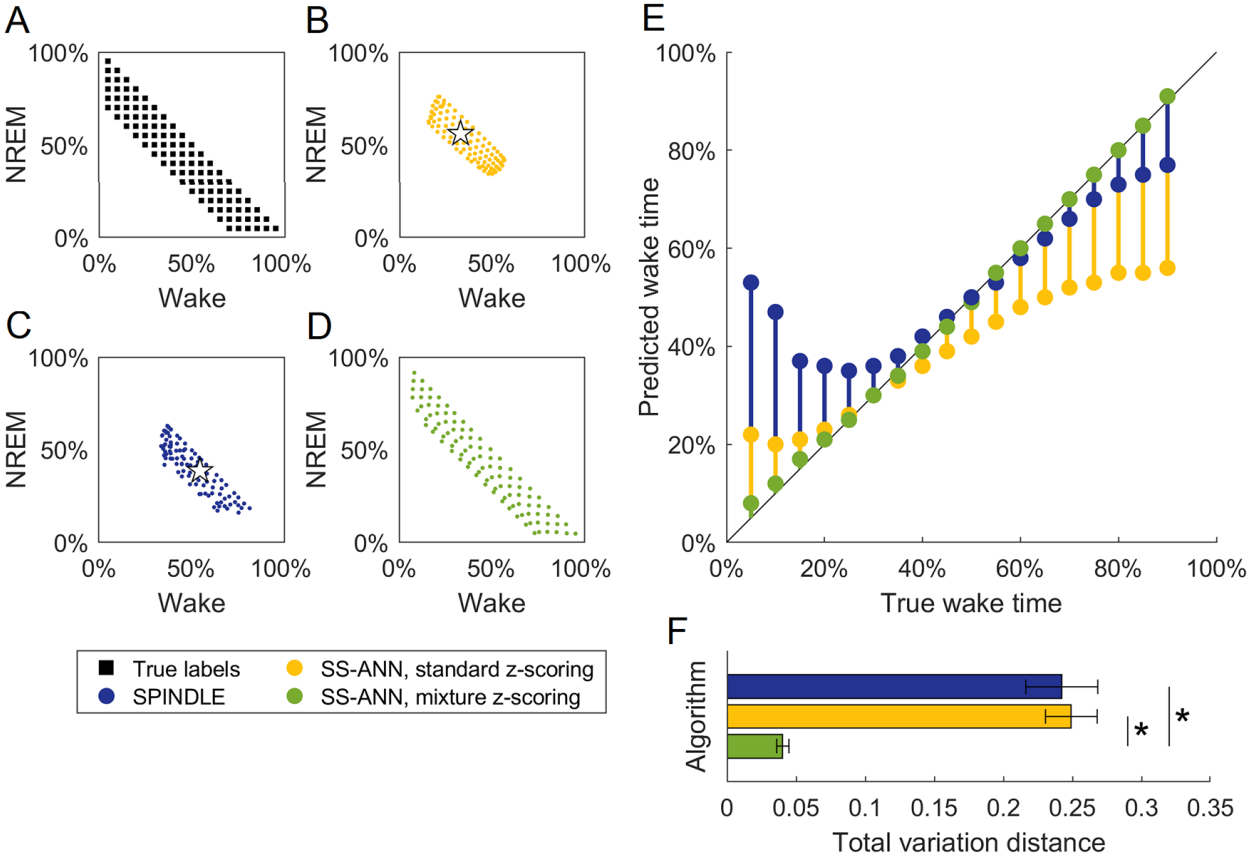


Figure 2.3. Comparison of sleep scoring algorithms on recordings with programmatically varied class balances.

Figure 2.4: **Validation of SS-ANN.** A: confusion matrix for SS-ANN on held-out data. The overall accuracy was 96.8%. The number of epochs is shown in parentheses. B: receiver operating characteristic (ROC) for SS-ANN. The inset panel shows a zoomed-in view of the upper-left corner. C: generalization across subjects. SS-ANN was trained on each mouse individually and tested on all mice. D: classification accuracy as a function of the amount of labeled data used for mixture  $z$ -scoring. Gray shading shows SEM,  $n = 20$  recordings.

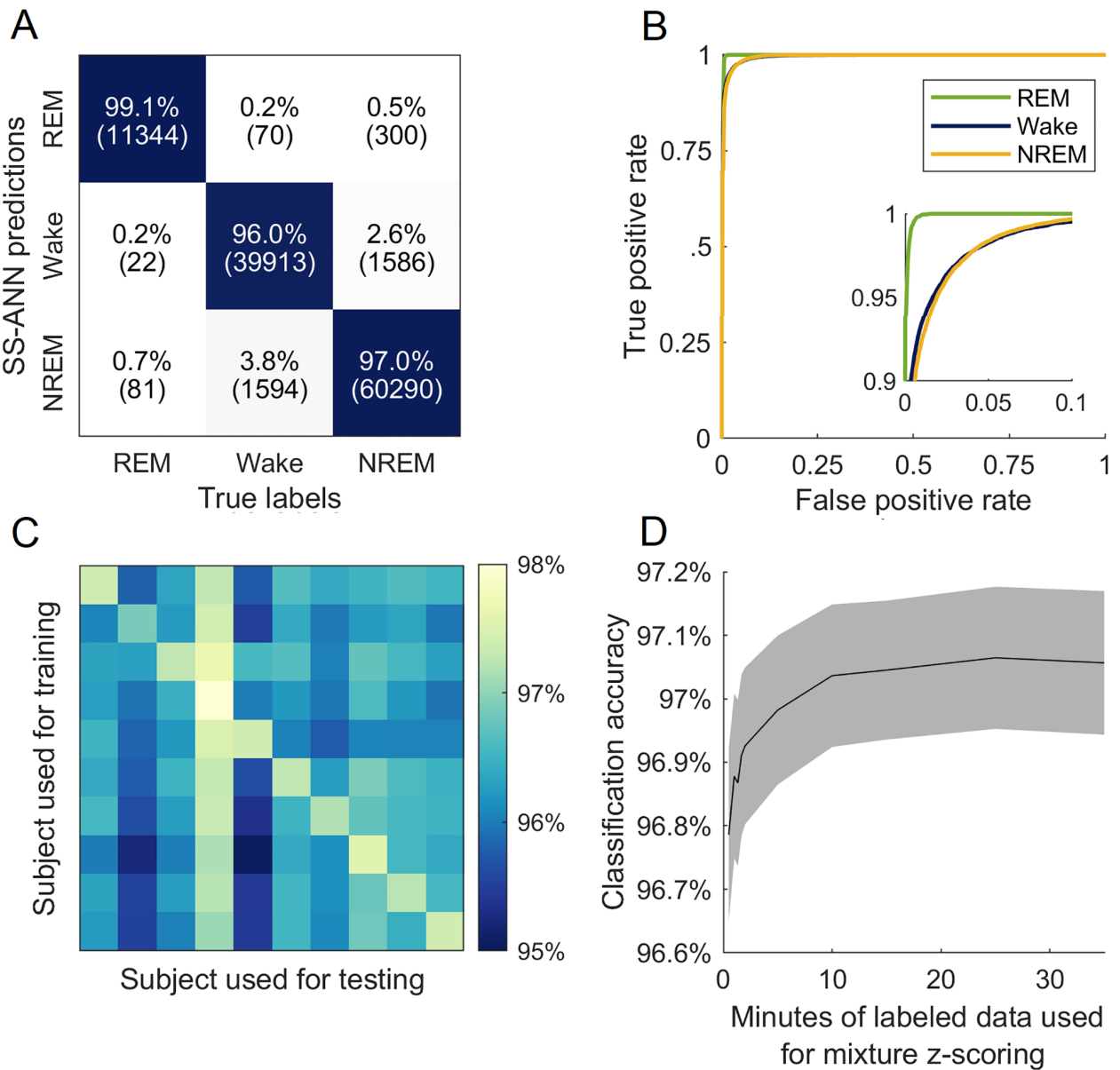


Figure 2.4. Validation of SS-ANN.

Figure 2.5: **AccuSleep interface for manual sleep scoring.** The lower three panels display the EEG and EMG signals as well as the sleep stage labels for epochs surrounding the currently selected epoch. The upper three panels provide context by displaying the sleep stages, EEG spectrogram, and EMG power on a longer time scale. The red line below the first panel indicates the time span of the lower three panels, and the diamond indicates the location of the currently selected epoch. For a complete description of the feature set of this software, please see the included user manual.

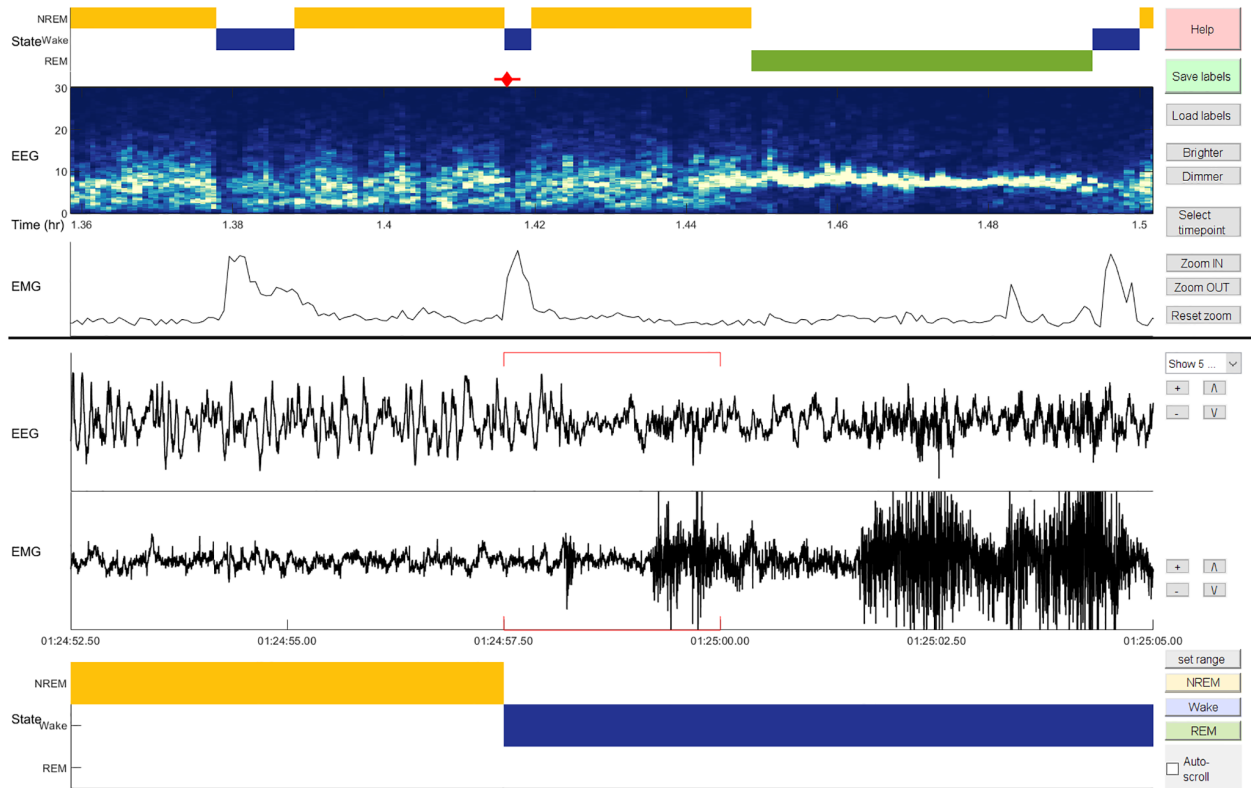


Figure 2.5. AccuSleep interface for manual sleep scoring.

Figure 2.6: **AccuSleep interface for automated sleep scoring.** After a small number of epochs are manually scored (Fig 2.5), the software uses SS-ANN and mixture  $z$ -scoring to perform automated sleep scoring on all recordings from a given subject simultaneously. Green check marks indicate valid user inputs.

**User manual**

Parameters for all recordings from one subject

Sampling rate (Hz): 512 ✓ Epoch length (sec): 2.5 ✓

**Recording list**

add remove

Recording 1  
Recording 2  
Recording 3  
Recording 4

Data / actions for the selected recording from this subject

Select EEG file D:\Data\Subject1\Day4\EEG.mat ✓

Select EMG file D:\Data\Subject1\Day4\EMG.mat ✓

Set / load label file D:\Data\Subject1\Day4\labels.mat ✓

Score selected manually ✓ ✓ ✓ ✓ ✓ Create calibration data file ✓ ✓ ✓ ✓ ?

Data / actions for all recordings from this subject

Load calibration data file D:\Data\Subject1\Day1\calibrationData.mat ✓

Load trained network file D:\Data\trained\_network.mat ✓

Score all automatically ✓ ✓ ✓ ✓ ✓ ✓ ✓  Only overwrite undefined epochs

Minimum bout length (sec): 5

**AccuSleep**

Messages

EMG file selected  
Label file found  
Inspecting EEG file...  
EEG file selected  
Inspecting EMG file...  
EMG file selected  
Label file found

Figure 2.6. AccuSleep interface for automated sleep scoring.



# Chapter 3

## Conclusion

### 3.1 Closing remarks

Inter-laboratory and inter-subject variability in PSG signals have limited the ability of automated sleep scoring algorithms to generalize. In this work, I presented a method, mixture  $z$ -scoring, that accounts for this variability and allows a compact neural network to perform highly accurate sleep scoring. Below, I address several remaining questions, and comment on the near-term outlook for sleep scoring.

Could the scoring algorithm be further refined? It is clear from the example training images in Fig 2.1 that some redundancy in the input feature space could be removed, especially in the higher frequency bands. However, even though we did not attempt to optimize the features used for classification or even the structure of the network itself, training the network is easily accomplished by a standard desktop computer (even without the use of a graphics processing unit) and classification of new data is likewise very fast—approximately 2-3 hours of PSG data can be scored per second. Since classification is therefore not a limiting step in the scoring process, and classification accuracy was at the level of inter-scorer agreement, we did not feel that further refinement was needed.

Could data from other species, such as humans or non-human primates, be scored using the same data processing pipeline? In subjects that show a mapping between sleep stages and the EEG spectrum that differs substantially from the mapping observed in mice, this would represent a case of concept shift—a class of dataset shift that cannot be corrected by mixture  $z$ -scoring. However, the AccuSleep software package includes a module that allows a user to train a new network using a set of labeled examples. Therefore, as long as the different sleep stages are distinguishable based on their EEG spectra and EMG patterns, classification should be possible. Given that scoring human sleep relies on brief EEG events that might not be detectable in spectrograms, and that spectrogram-based human sleep scoring by Vilamala et al. only attained 86% accuracy [1], it is unclear whether the algorithm presented

here would be suitable for this purpose.

What are the drawbacks of AccuSleep compared to other available scoring algorithms? As discussed in Chapter 2, mixture  $z$ -scoring requires a small amount of labeled data from each subject. While this step does require manual inspection of the PSG data, the cost is outweighed by a significant reduction in bias compared to methods that use standard  $z$ -scoring (Fig 1.3). In practice, this labeling step is likely to be problematic only for large-scale screens of sleep behavior [2–4]—although bias is clearly undesirable in these studies, as well.

## 3.2 Outlook

Now that convolutional neural networks operating on spectrograms of EEG and EMG data have been demonstrated by multiple laboratories to achieve accuracy equivalent to inter-rater reliability [5], it is unlikely that further gains can be made without incorporating additional sources of data. One possible avenue for progress is the use of video recordings. By processing videos of mice in their home cages with a CNN, Liu et al. were able to detect several sub-stages of wakefulness: locomotion, non-locomotor movement, and quiet wakefulness [6]. Reductions in the costs of graphics processing units and data storage, as well as the use of wireless recording devices, will make this approach more accessible. It may also be possible to use video data to measure respiration rate by magnifying small movements, at least during periods of quiescence [7].

Consistency across laboratories is another aspect of sleep scoring that could be improved. While algorithms for scoring human sleep are frequently benchmarked using the PhysioNet Sleep-EDF database [8], there is currently no equivalent database for rodent recordings. The PSG data collected as part of this work are freely available online at <https://osf.io/py5eb/> and could be used as a point of reference when evaluating the performance of different classifiers. Furthermore, because neural network models trained using AccuSleep can easily be shared and redeployed, the software presented here allows sleep researchers to share their expertise with laboratories anywhere in the world.

## References

- [1] Vilamala, A., Madsen, K. H., & Hansen, L. K. (2017). Deep convolutional neural networks for interpretable analysis of EEG sleep stage scoring. *arXiv e-prints*, arXiv 1710.00633.
- [2] Joshi, S. S., Sethi, M., Striz, M., Cole, N., Denegre, J. M., Ryan, J., Lhamon, M. E., Agarwal, A., Murray, S., Braun, R. E., Fardo, D. W., Kumar, V., Donohue, K. D., Sunderam, S., Chesler, E. J., Svenson, K. L., & O’Hara, B. F. (2019). Noninvasive sleep monitoring in large-scale screening of knock-out mice reveals novel sleep-related genes. *bioRxiv*. <https://doi.org/10.1101/517680>

- [3] Funato, H., Miyoshi, C., Fujiyama, T., Kanda, T., Sato, M., Wang, Z., Ma, J., Nakane, S., Tomita, J., Ikkyu, A. Et al. (2016). Forward-genetics analysis of sleep in randomly mutagenized mice. *Nature*, *539*(7629), 378–383.
- [4] Miyoshi, C., Kim, S. J., Ezaki, T., Ikkyu, A., Hotta-Hirashima, N., Kanno, S., Kakizaki, M., Yamada, M., Wakana, S., Yanagisawa, M., & Funato, H. (2019). Methodology and theoretical basis of forward genetic screening for sleep/wakefulness in mice, *116*(32), 16062–16067. <https://doi.org/10.1073/pnas.1906774116>
- [5] Miladinović, D., Muheim, C., Bauer, S., Spinnler, A., Noain, D., Bandarabadi, M., Gallusser, B., Krummenacher, G., Baumann, C., Adamantidis, A., Brown, S. A., & Buhmann, J. M. (2019). SPINDLE: End-to-end learning from EEG/EMG to extrapolate animal sleep scoring across experimental settings, labs and species. *PLOS Computational Biology*, *15*(4), 1–30. <https://doi.org/10.1371/journal.pcbi.1006968>
- [6] Liu, D., Li, W., Ma, C., Zheng, W., Yao, Y., Tso, C. F., Zhong, P., Chen, X., Song, J. H., Choi, W., Paik, S.-B., & Dan, Y. (2020). A common hub for sleep and motor control in the substantia nigra. *Science*, *367*(6476), 440–445.
- [7] Wu, H.-Y., Rubinstein, M., Shih, E., Guttag, J., Durand, F., & Freeman, W. (2012). Eulerian video magnification for revealing subtle changes in the world. *ACM transactions on graphics (TOG)*, *31*(4), 1–8.
- [8] 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.