



Modern Machine Learning as a Benchmark for Fitting Neural Responses

Ari S. Benjamin^{1*}, Hugo L. Fernandes², Tucker Tomlinson³, Pavan Ramkumar^{2,4}, Chris VerSteeg⁵, Raees H. Chowdhury^{3,5}, Lee E. Miller^{2,3,5} and Konrad P. Kording^{1,6}

¹ Department of Bioengineering, University of Pennsylvania, Philadelphia, PA, United States, ² Department of Physical Medicine and Rehabilitation, Rehabilitation Institute of Chicago, Northwestern University, Chicago, IL, United States, ³ Department of Physiology, Northwestern University, Chicago, IL, United States, ⁴ Department of Neurobiology, Northwestern University, Evanston, IL, United States, ⁵ Department of Biomedical Engineering, Northwestern University, Evanston, IL, United States, ⁶ Department of Neuroscience, University of Pennsylvania, Philadelphia, PA, United States

OPEN ACCESS

Edited by:

Yoram Burak,
Hebrew University of Jerusalem, Israel

Reviewed by:

Tatyana Sharpee,
Salk Institute for Biological Studies,
United States
Jonas Kubilius,
KU Leuven, Belgium and
Massachusetts Institute of
Technology, United States

*Correspondence:

Ari S. Benjamin
aarri@seas.upenn.edu

Received: 05 October 2017

Accepted: 29 June 2018

Published: 19 July 2018

Citation:

Benjamin AS, Fernandes HL, Tomlinson T, Ramkumar P, VerSteeg C, Chowdhury RH, Miller LE and Kording KP (2018) Modern Machine Learning as a Benchmark for Fitting Neural Responses. *Front. Comput. Neurosci.* 12:56. doi: 10.3389/fncom.2018.00056

Neuroscience has long focused on finding encoding models that effectively ask “what predicts neural spiking?” and generalized linear models (GLMs) are a typical approach. It is often unknown how much of explainable neural activity is captured, or missed, when fitting a model. Here we compared the predictive performance of simple models to three leading machine learning methods: feedforward neural networks, gradient boosted trees (using XGBoost), and stacked ensembles that combine the predictions of several methods. We predicted spike counts in macaque motor (M1) and somatosensory (S1) cortices from standard representations of reaching kinematics, and in rat hippocampal cells from open field location and orientation. Of these methods, XGBoost and the ensemble consistently produced more accurate spike rate predictions and were less sensitive to the preprocessing of features. These methods can thus be applied quickly to detect if feature sets relate to neural activity in a manner not captured by simpler methods. Encoding models built with a machine learning approach accurately predict spike rates and can offer meaningful benchmarks for simpler models.

Keywords: encoding models, neural coding, tuning curves, machine learning, generalized linear model, GLM, spike prediction

INTRODUCTION

A central tool of neuroscience is the tuning curve, which maps aspects of external stimuli to neural responses. The tuning curve can be used to determine what information a neuron encodes in its spikes. For a tuning curve to be meaningful it is important that it accurately describes the neural response. Often, however, methods are chosen for simplicity but not evaluated for their relative accuracy. Since inaccurate methods may systematically miss aspects of the neural response, any choice of predictive method should be compared with accurate benchmark methods.

A popular predictive model for neural data is the Generalized Linear Model (GLM) (Nelder and Baker, 1972; Simoncelli et al., 2004; Truccolo et al., 2005; Wu et al., 2006; Gerwinn et al., 2010). The GLM performs a nonlinear operation upon a linear combination of the input features, which are often called external covariates. Typical covariates are stimulus features, movement vectors, or the animal's location, and may include covariate history or spike history. In the absence of history terms, the GLM is also referred to as a linear-nonlinear Poisson (LN) cascade. The nonlinear

operation is usually held fixed, though it can be learned (Chichilnisky, 2001; Paninski et al., 2004a), and the linear weights of the combined inputs are chosen to maximize the agreement between the model fit and the neural recordings. This optimization problem of weight selection is convex, allowing a global optimum, and can be solved with efficient algorithms (Paninski, 2004). The assumption of Poisson firing statistics can often be loosened (Pillow et al., 2005), as well, allowing the modeling of a broad range of neural responses. Due to its ease of use, perceived interpretability, and flexibility, the GLM has become a popular model of neural spiking.

When using a GLM, it is important to check that the method's assumptions about the data are correct. The GLM's central assumption is that the inputs relate linearly to the log firing rate, or generally some monotonic function of the firing rate. It thus cannot learn arbitrary multi-dimensional functions of the inputs. When the nonlinearity is different than assumed, it is likely that the optimal weight on one input will depend on the values of other inputs. In this case the GLM will only partially represent the neural response, will poorly predict activity, and may not be reproducible on other datasets. This drawback has been noted before, and indeed the GLM has been shown to miss nonlinearity in numerous circumstances (Butts et al., 2011; Freeman et al., 2015; Heitman et al., 2016; McIntosh et al., 2016). However, GLMs are still commonly applied without comparison to other methods. To test if the linearity assumption is valid, it is sufficient to test if other nonlinear methods predict activity more accurately from the same features. Many extensions have been proposed that introduce a specific form of nonlinearity (McFarland et al., 2013; Theis et al., 2013; Latimer et al., 2014; Williamson et al., 2015; Maheswaranathan et al., 2017), but these methods ask specific research questions and are not intended as general benchmarks. What is needed is are nonlinear methods that are universally applicable to new data.

Machine learning (ML) methods for regression have improved dramatically since the invention of the GLM. Many ML methods require little feature engineering (i.e., pre-transformations the features) and do not need to assume linearity. These methods are thus ideal candidates for benchmark methods. The ML approach is now quite standardized and robust across many domains of data. As exemplified by winning solutions on Kaggle, an ML competition website (Kaggle Winner's Blog, 2016), the usual approach is to fit several top performing methods, and then to ensemble these models together. These methods are now relatively easy to implement in a few lines of code in a scripting language such as Python, and are enabled by well-supported machine learning packages, such as scikit-learn (Pedregosa et al., 2011), Keras (Chollet, 2015), Tensorflow (Abadi et al., 2016), and XGBoost (Chen and Guestrin, 2016). The greatly increased predictive power of modern ML methods is now very accessible and could help to benchmark and improve the state of the art in encoding models across neuroscience.

In order to investigate the feasibility of ML as a benchmark approach, we applied several ML methods, including artificial neural networks, gradient boosted trees, and ensembles to the task of predicting spike rates, and evaluated their performance alongside a GLM. We compared the methods on data from

three separate brain areas. These areas differed greatly in the effect size of covariates and in their typical spike rates, and thus served to evaluate the strengths of these methods across different conditions. In each area we found that the ensemble of methods could more accurately predict spiking than the GLM with typical feature choices. The use of an ML benchmark thus made clear that tuning curves built for these features with a GLM would not capture the full nature of neural activity. We provide our implementing code at <https://github.com/KordingLab/spykesML> so that all neuroscientists may easily test and compare ML to their own methods on other datasets.

MATERIALS AND METHODS

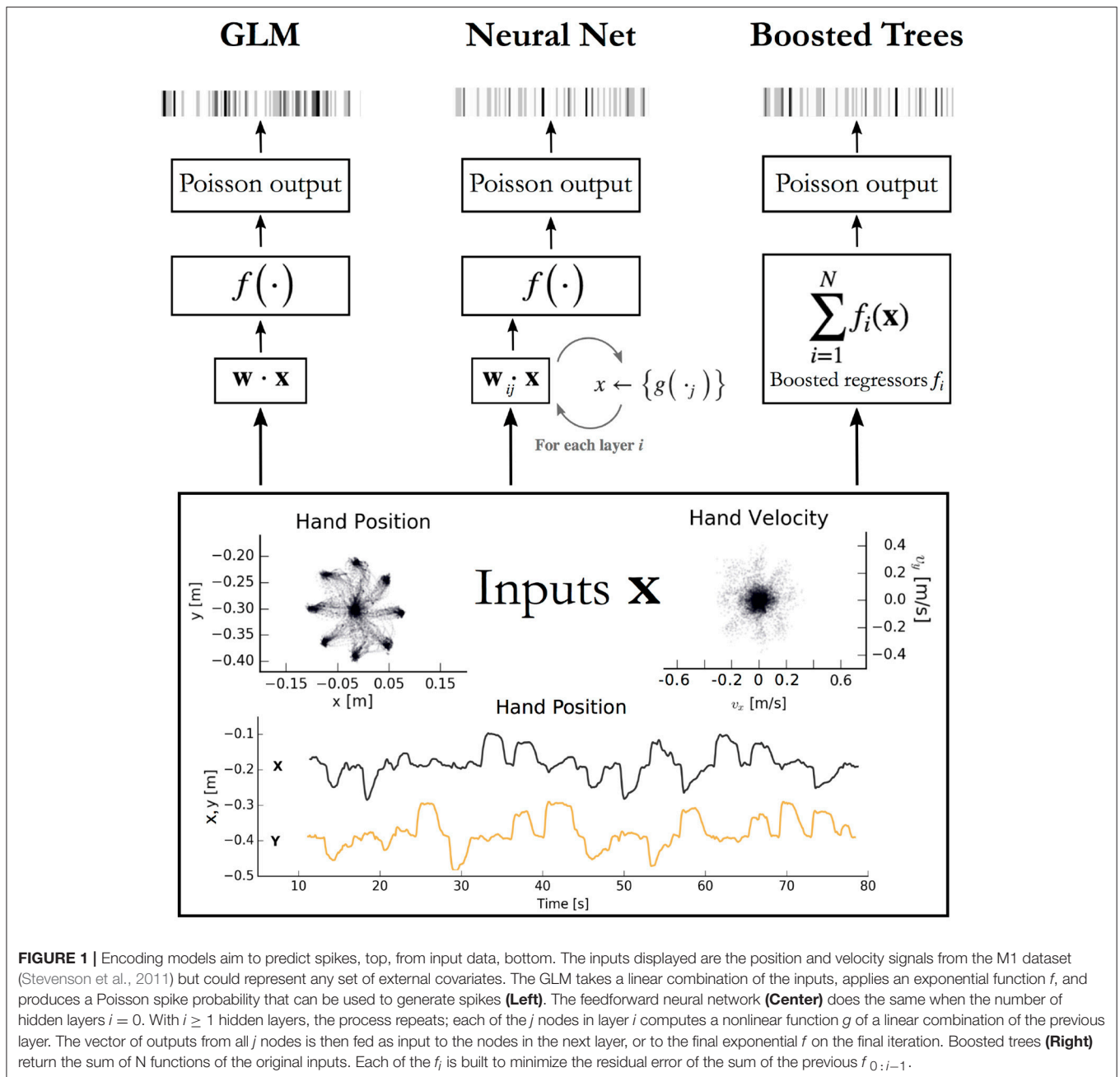
Data

We tested our methods at predicting spike rates for neurons in the macaque primary motor cortex, the macaque primary somatosensory cortex, and the rat hippocampus. All animal use procedures were approved by the institutional animal care and use committees at Northwestern University and conform to the principles outlined in the Guide for the Care and Use of Laboratory Animals (National Institutes of Health publication no. 86-23, revised 1985). Data presented here were previously recorded for use with multiple analyses. Procedures were designed to minimize animal suffering and reduce the number used.

The macaque motor cortex data consisted of previously published electrophysiological recordings from 82 neurons in the primary motor cortex (M1) (Stevenson et al., 2011). The neurons were sorted from recordings made during a two-dimensional center-out reaching task with eight targets. In this task the monkey grasped the handle of a planar manipulandum that controlled a cursor on a computer screen and simultaneously measured the hand location and velocity (Figure 1). After training, an electrode array was implanted in the arm area of area 4 on the precentral gyrus. Spikes were discriminated using offline sorter (Plexon, Inc), counted and collected in 50-ms bins. The neural recordings used here were taken in a single session lasting around 13 min.

The macaque primary somatosensory cortex (S1) data was recorded during a two-dimensional random-pursuit reaching task and was previously unpublished. In this task, the monkey gripped the handle of the same manipulandum. The monkey was rewarded for bringing the cursor to a series of randomly positioned targets appearing on the screen. After training, an electrode array was implanted in the arm area of area 2 on the post-central gyrus, which receives a mix of cutaneous and proprioceptive afferents. Spikes were processed as for M1. The data used for this publication derives from a single recording session lasting 51 min.

As with M1 (described in results), we processed the hand position, velocity, and acceleration accompanying the S1 recordings in an attempt to obtain linearized features. The features (x, y, \dot{x}, \dot{y}) were found to be the most successful for the GLM. Since cells in the arm area of S1 have been shown to have approximately sinusoidal tuning curves relating to movement direction (Prud'homme and Kalaska, 1994), we also tested the



same feature transformations as were performed for M1 but did not observe any increase in predictive power.

The third dataset consists of recordings from 58 neurons in the CA1 region of the rat dorsal hippocampus during a single 93 min free foraging experiment, previously published and made available online by the authors (Mizuseki et al., 2009a,b). Position data from two head-mounted LEDs provided position and heading direction inputs. Here we binned inputs and spikes from this experiment into 50 ms bins. Since many neurons in the dorsal hippocampus are responsive to the location of the rat, we processed the 2D position data into a list of squared distances from a 5×5 grid of place fields that tile the workspace. Each

position feature thus has the form

$$p_{ij} = \frac{1}{2} (x(t) - \mu_{ij})^T \Sigma_{ij}^{-1} (x(t) - \mu_{ij}),$$

where μ_{ij} is the center of place field i , $j \leq 5$ and Σ_{ij} is a covariance matrix chosen for the uniformity of tiling. An exponentiated linear combination of the p_{ij} (as is performed in the GLM) evaluates to a single Gaussian centered anywhere between the place fields. The inclusion of the p_{ij} as features thus transforms the standard representation of cell-specific place fields (Brown et al., 1998) into the mathematical formulation of a GLM. The

final set of features included the p_{ij} as well as the rat speed and head orientation.

Treatment of Spike and Covariate History

We slightly modified our data preparation methods for spike rate prediction when spike and covariate history terms were included as regressors (**Figure 6**). To construct spike and covariate history filters, we convolved 10 raised cosine bases (built as in Pillow et al., 2008) with binned spikes and covariates. The longest temporal basis included times up to 250 ms before the time bin being predicted. This process resulted in 120 total covariates per sample (10 current covariates, 100 covariate temporal filters, and 10 spike history filters). We predicted spike rates in 5 ms bins (rather than 50 ms) to allow for modeling of more precise time-dependent phenomena, such as refractory effects. The cross-validation scheme also differs from the main analysis of this paper, as using randomly selected splits of the data would result in the appearance in the test set of samples that were in history terms of training sets, potentially resulting in overfitting. We thus employed a cross-validation routine to split the data continuously in time, assuring that no test set sample has appeared in any form in training sets.

Generalized Linear Model

The Poisson GLM is a multivariate regression model that describes the instantaneous firing rate as a nonlinear function of a linear combination of input features (see e.g., Schwartz et al., 2006; Aljadeff et al., 2016 for review, Pillow et al., 2008; Fernandes et al., 2014; Ramkumar et al., 2016 for usage). Here, we took the form of the nonlinearity to be exponential, as is common in previous applications of GLMs to similar data (Saleh et al., 2012). It should be noted that it is also possible to learn arbitrary link functions through histogram methods (Chichilnisky, 2001; Paninski et al., 2004a). We approximate neural activity as a Poisson process, in which the probability of firing in any instant is independent of firing history. The general form of the GLM is depicted in **Figure 1**. We implemented the GLM using elastic-net regularization, using the open-source Python package `pyglmnet` (Ramkumar et al., 2017). The regularization path was optimized separately on a single neuron in each dataset on a validation set not used for scoring.

Neural Network

Neural networks are well-known for their success at supervised learning tasks. More comprehensive reviews can be found elsewhere (Schmidhuber, 2015). Here, we implemented a simple feedforward neural network and, for the analysis with history terms, an LSTM, a recurrent neural network architecture that allows the modeling of time dependencies on multiple time-scales (Gers et al., 2000).

We point out that a feedforward neural network with no hidden layers is equivalent in mathematical form to a GLM (**Figure 1**). For multilayer networks, one can write each hidden layer of n nodes as simply n GLMs, each taking the output of the previous layer as inputs (noting that the weights of each are chosen to maximize only the final objective function, and that the intermediate nonlinearities need not be the same as the output

nonlinearity). A feedforward neural network can be seen as a generalization, or repeated application of a GLM.

The networks were implemented with the open-source neural network library Keras, running Theano as the backend (Chollet, 2015; Team et al., 2016). The feedforward network contained two hidden layers, dense connections, rectified linear activation, and a final exponentiation. To help avoid overfitting, we allowed dropout on the first layer, included batch normalization, and allowed elastic-net regularization upon the weights (but not the bias term) of the network (Srivastava et al., 2014). The networks were trained to maximize the Poisson likelihood of the neural response. We optimized over the number of nodes in the first and second hidden layers, the dropout rate, and the regularization parameters for the feedforward neural network, and for the number of epochs, units, dropout rate, and batch size for the LSTM. Optimization was performed on only a subset of the data from a single neuron in each dataset, using Bayesian optimization (Snoek et al., 2012) in an open-source Python implementation (BayesianOptimization, 2016).

Gradient Boosted Trees

A popular method in many machine learning competitions is that of gradient boosted trees. Here we describe the general operation of XGBoost, an open-source implementation that is efficient and highly scalable, works on sparse data, and easy to implement out-of-the-box (Chen and Guestrin, 2016).

XGBoost trains many sequential models to minimize the residual error of the sum of previous model. Each model is a decision tree, or more specifically a classification and regression tree (CART) (Friedman, 2001). Training a decision tree amounts to determining a series of rule-based splits on the input to classify output. The CART algorithm generalizes this to regression by taking continuously-valued weights on each of the leaves of the decision tree.

For any predictive model $\hat{y}^{(1)} = f_1(\mathbf{x}_i)$ and true response y_i , we can define a loss function $l(\hat{y}^{(1)}, y_i)$ between the prediction and the response. The objective to be minimized during training is then simply the sum of the loss over each training example i , plus some regularizing function Ω that biases toward simple models.

$$L = \sum_i l(\hat{y}_i^{(1)}, y_i) + \Omega(f_1)$$

After minimizing L for a single tree, XGBoost constructs a second tree $f_2(\mathbf{x}_i)$ that approximates the residual. The objective to be minimized is thus the total loss L between the true response y_i and the sum of the predictions given by the first tree and the one to be trained.

$$L = \sum_i l(\hat{y}_i^{(1)} + f_2(\mathbf{x}_i), y_i) + \Omega(f_2)$$

This process is continued sequentially for a predetermined number of trees, each trained to approximate the residual of the sum of previous trees. In this manner XGBoost is designed to progressively decrease the total loss with each additional tree. At

the end of training, new predictions are given by the sum of the outputs of all trees.

$$\hat{y} = \sum_{k=1}^N f_k(\mathbf{x})$$

In practice, it is simpler to choose the functions f_k via gradient boosting, which minimizes a second order approximation of the loss function (Friedman et al., 2000).

XGBoost offers several additional parameters to optimize performance and prevent overfitting. Many of these describe the training criteria for each tree. We optimized some of these parameters for a single neuron in each dataset using Bayesian optimization (again over a validation set different from the final test set). These parameters included the number of trees to train, the maximum depth of each decision tree, and the minimum weight allowed on each decision leaf, the data subsampling ratio, and the minimum gain required to create a new decision branch.

Random Forests

We implement random forests here to increase the power of the ensemble (see below); their performance alone is displayed in Supplementary Figure 1. It should be noted that the Scikit-learn implementation currently only minimizes the mean-squared error of the output, which is not properly applicable to Poisson processes and may cause poor performance. Despite this drawback their presence still improves the ensemble scores. Random forests train multiple parallel decision trees on the features-to-spikes regression problem (not sequentially on the remaining residual, as in XGBoost) and averages their outputs (Ho, 1998). The variance on each decision tree is increased by training on a sample of the data drawn with replacement (i.e., bootstrapped inputs) and by choosing new splits using only a random subset of the available features. Random forests are implemented in Scikit-learn (Pedregosa et al., 2011).

Ensemble Method

It is a common machine learning practice to create ensembles of several trained models. Different algorithms may learn different characteristics of the data, make different types of errors, or generalize differently to new examples. Ensemble methods allow for the successes of different algorithms to be combined. Here we implemented *stacking*, in which the output of several models is taken as the input set of a new model (Wolpert, 1992). After training the GLM, neural network, random forest, and XGBoost on the features of each dataset, we trained an additional instance of XGBoost using the spike rate predictions of the previous methods as input. The outputs of this “second stage” XGBoost are the predictions of the ensemble.

Scoring and Cross-Validation

Each of the three methods was scored with the Poisson pseudo- R^2 score, a scoring function applicable to Poisson processes (Cameron and Windmeijer, 1997). Note that a standard R^2 score

assumes Gaussian noise and cannot be applied here. The pseudo- R^2 was calculated as one minus the ratio of the deviances of the predicted output \hat{y} to the mean firing rate \bar{y} .

$$R_M^2 = 1 - \frac{D(\hat{y})}{D(\bar{y})}$$

We can gain intuition into the pseudo- R^2 score by writing out the deviances in terms of log likelihoods $L()$, and combining the fraction.

$$R_M^2 = 1 - \frac{\log L(y) - \log L(\hat{y})}{\log L(y) - \log L(\bar{y})} = \frac{\log L(\hat{y}) - \log L(\bar{y})}{\log L(y) - \log L(\bar{y})}$$

This expression includes $L(y)$, which is the log likelihood of the “saturated model,” which offers one parameter per observation and models the data perfectly. The pseudo- R^2 can thus be interpreted as the fraction of the maximum potential log-likelihood gain achieved by the tested model (Cameron and Windmeijer, 1997). It takes a value of 0 when the data is as likely under the tested model as the null model, and a value of 1 when the tested model perfectly describes the data. It is empirically a lower value than a standard R^2 when both are applicable (Domencich and McFadden, 1975). The null model can also be taken to be a model other than the mean firing rate (e.g., the GLM) to directly compare two methods, in which case we refer to the score as the “comparative pseudo- R^2 .” The comparative pseudo- R^2 is referred to elsewhere as the “relative pseudo- R^2 ,” renamed here to avoid confusion with the difference of two standard pseudo- R^2 scores both measured against the mean (Fernandes et al., 2014).

We used 8-fold cross-validation (CV) when assigning a final score to the models. The input and spike data were segmented into eight equal partitions. These partitions were continuous in time when spike and covariate history were included as covariates, and otherwise were segmented randomly in time. The methods were trained on seven partitions and tested on the eighth, and this was repeated until all segments served as the test partition once. The mean of the eight scores are then recorded for the final score.

Cross-validation for ensemble methods requires extra care since the inputs for the ensemble are themselves model predictions for each data point. The training set for the ensemble must contain predictions from methods that were themselves not trained on the validation set. Otherwise, there may be a leak of information from the validation set into the training set and the validation score might be better than on a true held-out set. This rules out using simple k -fold CV with all methods and the ensemble trained on the same test/train splits. Instead, we used a nested CV scheme to train and score the ensemble. We create an outer $j = 8$ folds to build training and test sets for the ensemble. On each outer fold we create first-order predictions for each data point in the following manner. We first run an inner k -fold CV on just the training set (i.e., 7/8 of the original dataset) with each first stage method such that we obtain predictions for the whole

training set of that fold. This ensures that the ensemble's test set was never used for training any method. Finally, we build the ensemble's test set from the predictions of the first stage methods trained on the entire training set. The ensemble can then be tested on a held-out set that was never used to fit any model. The process is repeated for each of the j folds and the mean and variance of the j scores of the ensemble's predictions are recorded.

RESULTS

We applied several machine learning methods to predict spike counts in three brain regions and compared the quality of the predictions to those of a GLM. Our primary analysis centered on neural recordings from the macaque primary motor cortex (M1) during reaching (**Figure 1**). We examined the methods' relative performance on several sets of movement features with various levels of preprocessing, including one set that included spike and covariate history terms. Analyses of data from rhesus macaque S1 and rat hippocampus indicate how these methods compare for areas other than M1. On each of the three datasets we trained a GLM and compared it to the performance of a feedforward neural network, XGBoost (a gradient boosted trees implementation), and an ensemble method. The ensemble was an additional instance of XGBoost trained on the predictions of all three methods plus a random forest regressor. The application of these methods allowed us to demonstrate the potential of a modern approach to be able to identify whether there are typically neural nonlinearities that are not captured by a GLM. The code implementing these methods can be used by any electrophysiology lab to benchmark their own encoding models.

To test that all methods work reasonably well in a trivial case, we trained each to predict spiking from a simple, well-understood feature. Some neurons in M1 have been described as responding linearly to the exponentiated cosine of movement direction relative to a preferred angle (Amirikian and Georgopoulos, 2000). We therefore predicted the spiking of M1 neurons from the cosine and sine of the direction of hand movement in the reaching task. (The linear combination of a sine and cosine curve is a phase-shifted cosine, by identity, allowing the GLM to learn the proper preferred direction). We observed that each method identified a similar tuning curve (**Figure 2B**) and that the bulk of the neurons in the dataset were just as well predicted by each of the methods (**Figures 2A,C**) {though the ensemble was slightly more accurate than the GLM, with mean comparative pseudo- R^2 above zero, 0.06 [0.043 – 0.084], 95% bootstrapped confidence interval (CI)}. The similar performance suggested that, for the majority of neurons, an exponentiated cosine successfully approximates the response to movement direction alone, as has been previously found (Paninski et al., 2004b). All methods can in principle estimate tuning curves, and machine learning can indicate if the proper features are used.

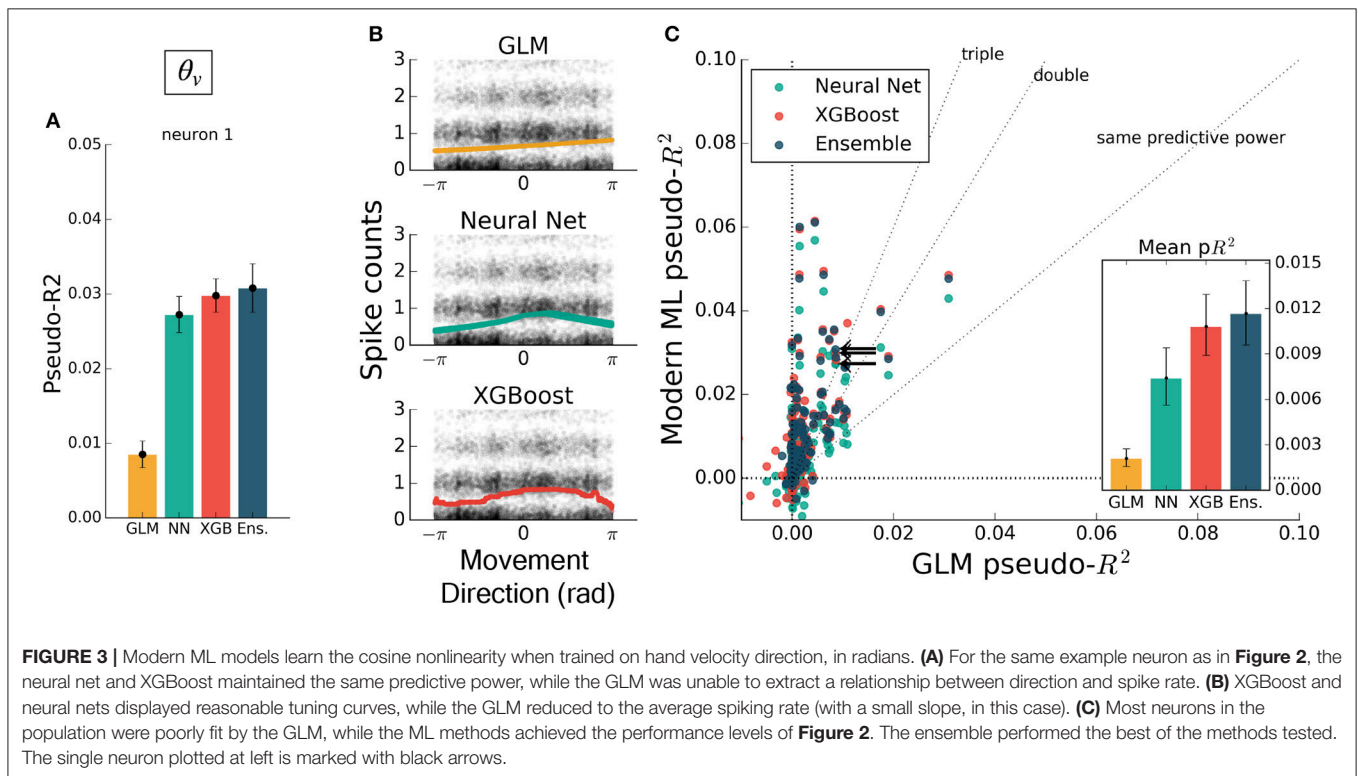
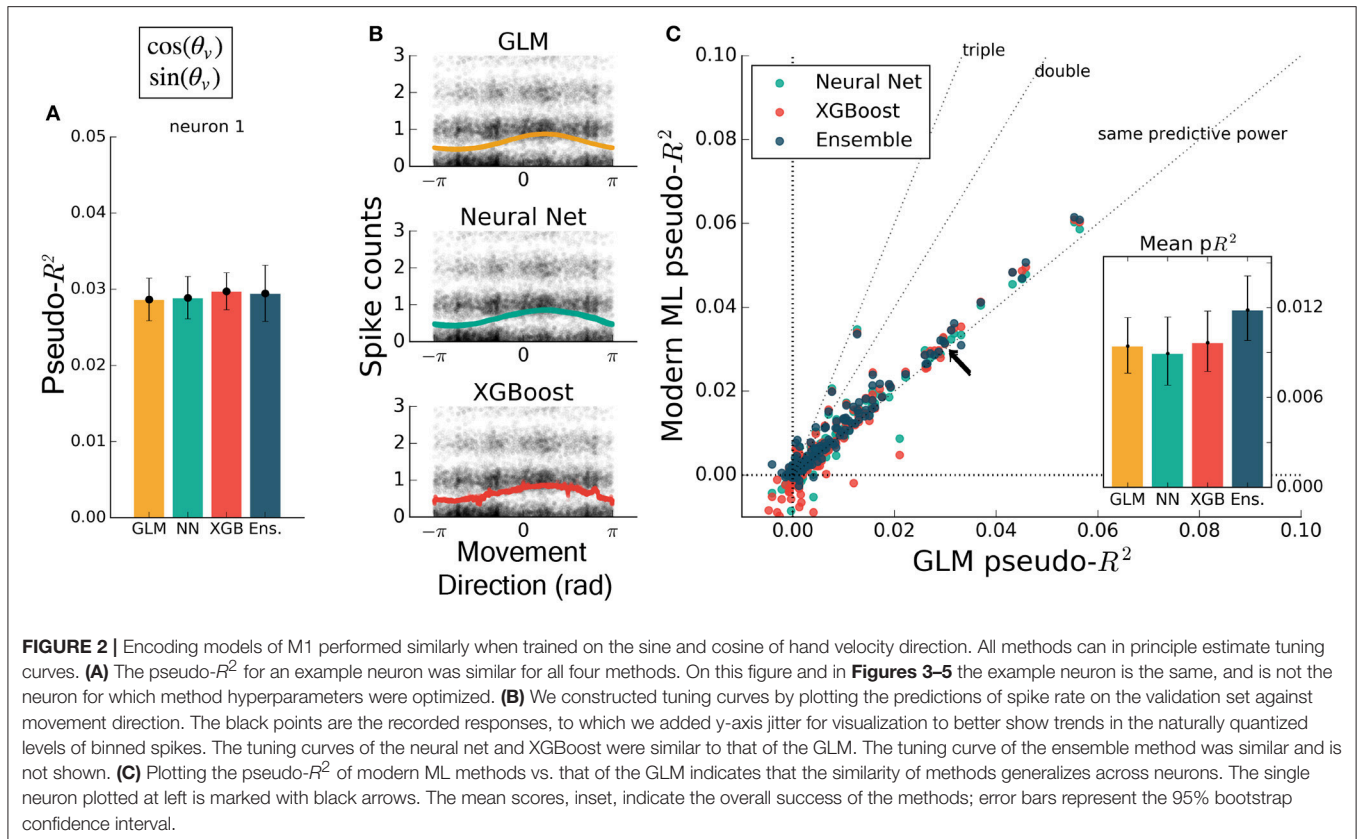
If the form of the nonlinearity is not known, machine learning can still attain good predictive ability. To illustrate the ability of modern machine learning to find the proper nonlinearity, we performed the same analysis as above but omitted the initial cosine feature-engineering step. Trained on only the hand

velocity direction, in radians, which changes discontinuously at $\pm\pi$, all methods but the GLM closely matched the predictive power they attained using the engineered feature (**Figure 3A**). The GLM failed at generating a meaningful tuning curve, which was expected since the exponentiated velocity direction is not equal to cosine tuning (**Figure 3B**). Both trends were consistent across the population of recorded neurons (**Figure 3C**). The neural net, XGBoost, and ensemble methods can learn the nonlinearity of single features without requiring manual feature transformation.

The inclusion of multiple features raises the possibility of nonlinear feature interactions that may elude a GLM. As a simple demonstration of this principle, we trained all methods on the four-dimensional set of hand position and velocity (x, y, \dot{x}, \dot{y}) . While all methods gained predictive power relative to models using movement direction alone, the GLM failed to match the other methods (**Figures 4A,C**). If the GLM was fit alone, and no further featuring engineering been attempted, these features would have appeared to be relatively uninformative of the neural response. If nonlinear interactions exist between preselected features, machine learning methods can potentially learn these interactions and indicate if more linearly-related features exist.

While feature engineering can improve the performance of GLMs, it is not always simple to guess the optimal set of processed features. We demonstrated this by training all methods on features that have previously been successful at explaining spike rate in a similar center-out reaching task (Paninski et al., 2004a). These extra features included the sine and cosine of velocity direction (as in **Figure 2**), and the speed, radial distance of hand position, and the sine and cosine of position direction. The training set was thus 10-dimensional, though highly redundant, and was aimed at maximizing the predictive power of the GLM. Feature engineering improved the predictive power of all methods to variable degrees, with the GLM improving to the level of the neural network (**Figure 5**). XGBoost and the ensemble still predicted spike rates better than the GLM (**Figure 5C**), with the ensemble scoring on average nearly double the GLM (ratio of population means of 1.8 [1.4 – 2.2], 95% bootstrapped CI). The ensemble was significantly better than XGBoost (mean comparative pseudo- R^2 of 0.08 [0.055 – 0.103], 95% bootstrapped CI) and was thus consistently the best predictor. Though standard feature engineering greatly improved the GLM, the ensemble and XGBoost still could identify that neural nonlinearity was missed by the GLM.

It is important to note that the specific ordering of methods depends on features such as the amount of data available for training. We investigated this dependence for the M1 dataset by plotting the cross-validated performance as a function of the fraction of the data used for training (Supplementary Figure 3). Some neurons are best fit by the GLM when very little data is available, while other neurons are best fit by XGBoost and the ensemble for any amount of data tested. The neural network is most sensitive to training data availability. This sensitivity to the domain of data emphasizes the importance of the applied ML paradigm of evaluating (and potentially ensembling) many methods.



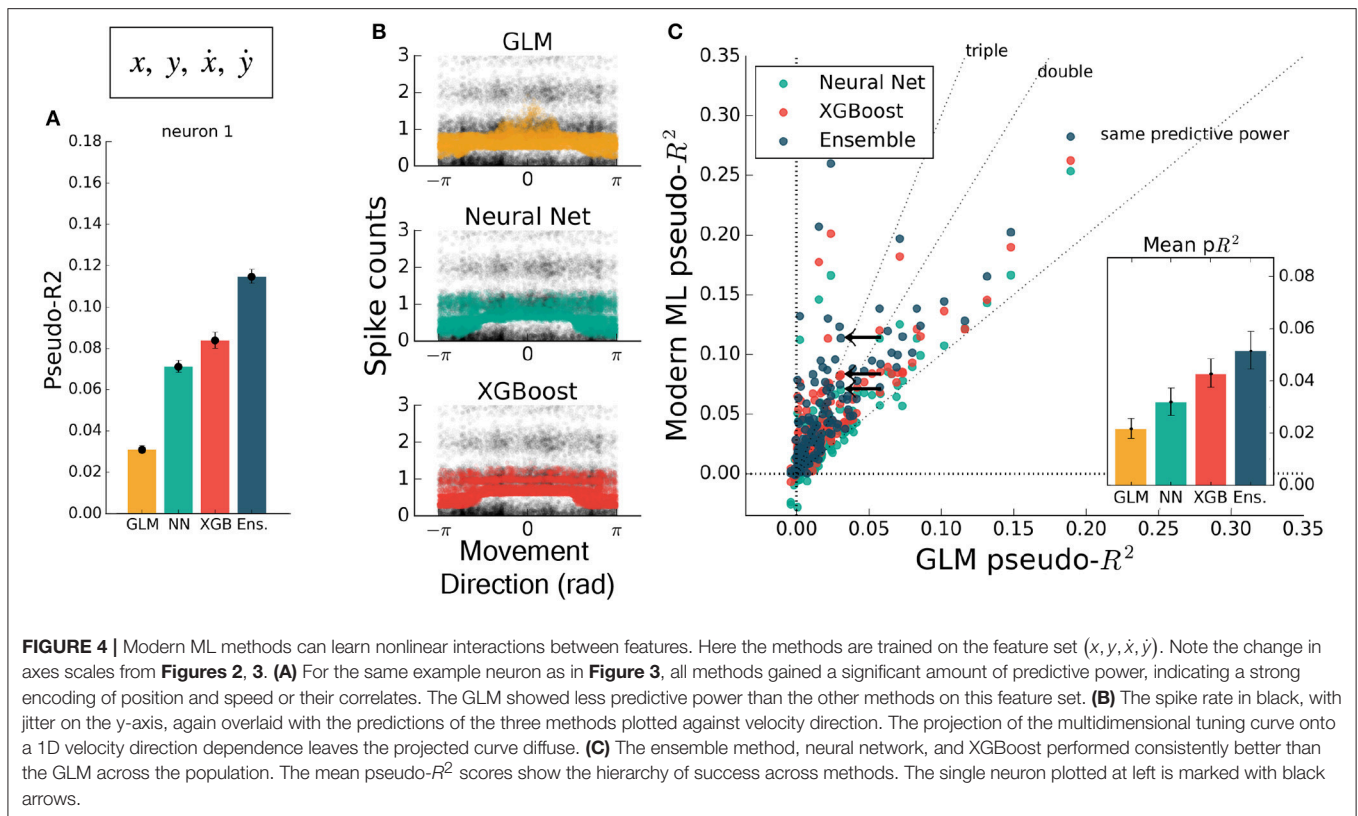


FIGURE 4 | Modern ML methods can learn nonlinear interactions between features. Here the methods are trained on the feature set (x, y, \dot{x}, \dot{y}) . Note the change in axes scales from **Figures 2, 3**. **(A)** For the same example neuron as in **Figure 3**, all methods gained a significant amount of predictive power, indicating a strong encoding of position and speed or their correlates. The GLM showed less predictive power than the other methods on this feature set. **(B)** The spike rate in black, with jitter on the y-axis, again overlaid with the predictions of the three methods plotted against velocity direction. The projection of the multidimensional tuning curve onto a 1D velocity direction dependence leaves the projected curve diffuse. **(C)** The ensemble method, neural network, and XGBoost performed consistently better than the GLM across the population. The mean pseudo- R^2 scores show the hierarchy of success across methods. The single neuron plotted at left is marked with black arrows.

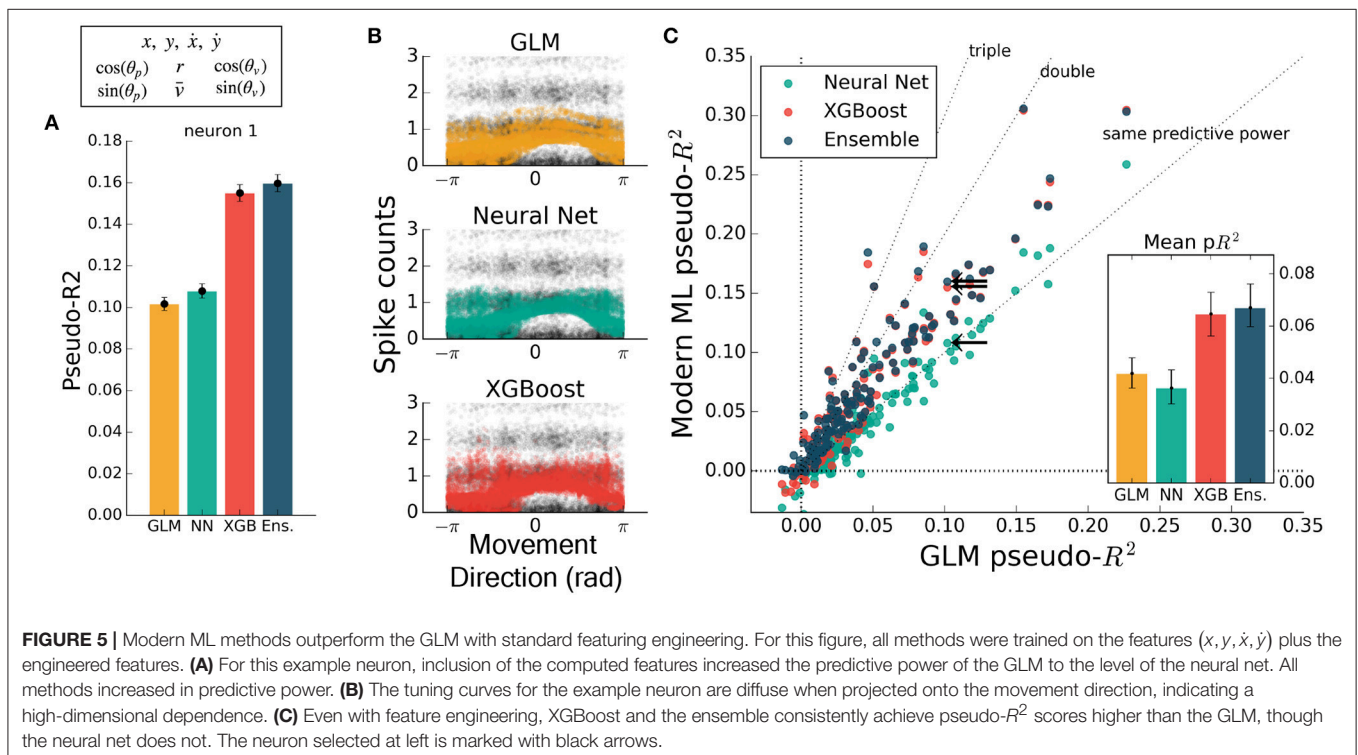
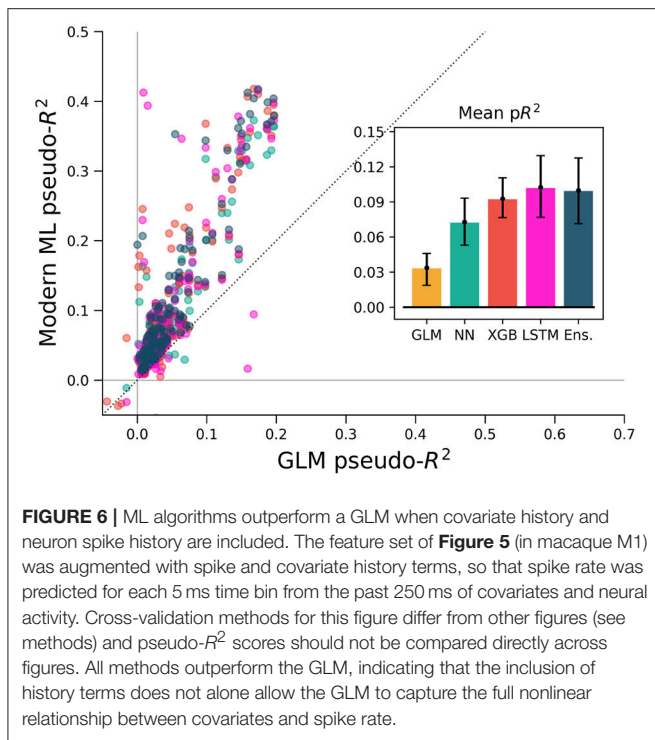


FIGURE 5 | Modern ML methods outperform the GLM with standard featuring engineering. For this figure, all methods were trained on the features (x, y, \dot{x}, \dot{y}) plus the engineered features. **(A)** For this example neuron, inclusion of the computed features increased the predictive power of the GLM to the level of the neural net. All methods increased in predictive power. **(B)** The tuning curves for the example neuron are diffuse when projected onto the movement direction, indicating a high-dimensional dependence. **(C)** Even with feature engineering, XGBoost and the ensemble consistently achieve pseudo- R^2 scores higher than the GLM, though the neural net does not. The neuron selected at left is marked with black arrows.

Studies employing a GLM often include activity history as a covariate when predicting spike rates, as well as past values

of the covariates themselves, and it is known that this allows GLMs to model a wider range of phenomena (Weber and Pillow,



2016). We tested various ML methods on the M1 dataset using this history-augmented feature set to see if all methods would still explain a similar level of activity. We binned data by 5 ms (rather than 50 ms) to agree in timescale with similar studies, and built temporal filters by convolving 10 raised-cosine bases with features and spikes. We note that smaller time bins result in a sparser dataset, and thus pseudo- R^2 scores cannot be directly compared with other analysis in this paper. On this problem, our selected ML algorithms again outperformed the GLM (**Figure 6**). The overall best algorithm was the LSTM, which we include here as it specifically designed for modeling time series, though for most neurons XGBoost performed similarly. Thus, for M1 neurons, the GLM did not capture all predictable phenomena even when spike and covariate history were included.

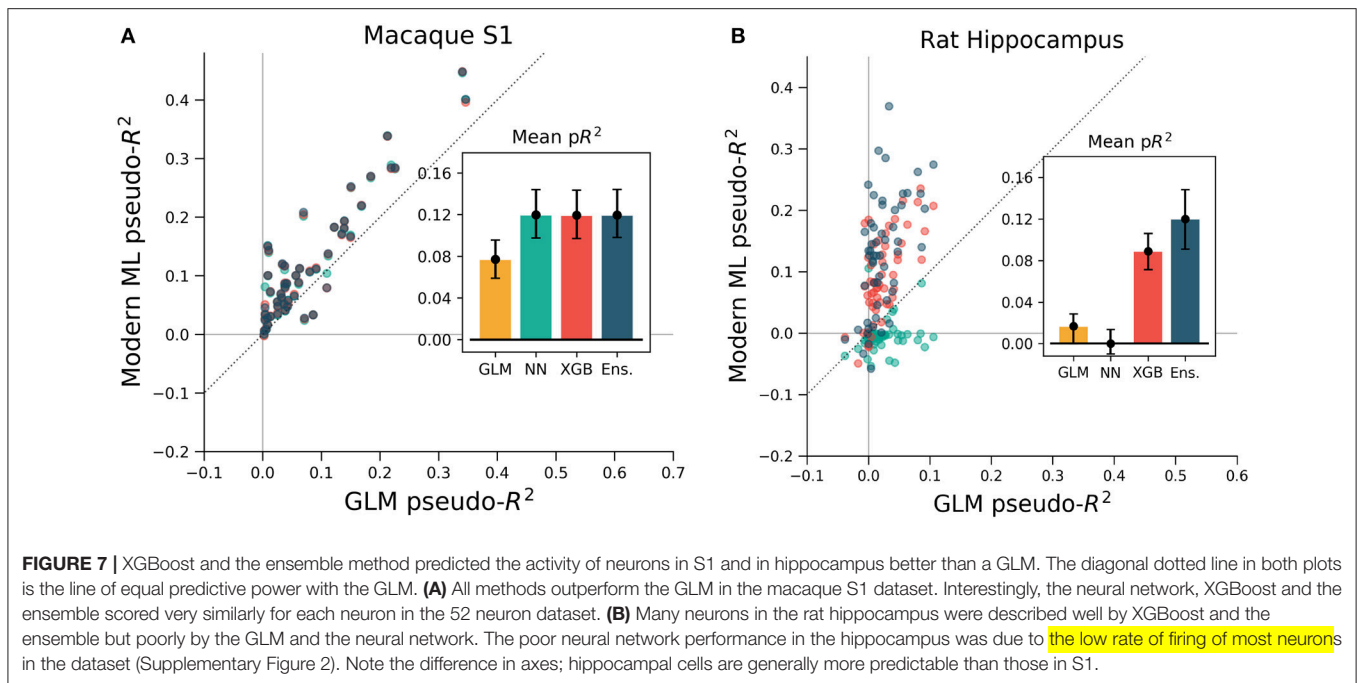
To ensure that these results were not specific to the motor cortex, we extended the same analyses to primary somatosensory cortex (S1) data. We again predicted neural activity from hand movement and speed, and here without spike or covariate history terms. The ML methods outperformed the GLM for all but three of the 52 neurons, indicating that firing rates in S1 generally relate nonlinearly to hand position and velocity (**Figure 7A**). Each of the three ML methods performed similarly for each neuron. The S1 neural function was thus equally learnable by each method, which is surprising given the dissimilarity of the neural network and XGBoost algorithms. This situation would occur if learning has saturated near ground truth, though this cannot be proven definitively to be the case. It is at least clear from the underperformance of the GLM that the relationship of S1 activity to these covariates is nonlinear beyond the assumptions of the GLM.

We asked if the same trends of performance would hold for the rat hippocampus dataset, which was characterized by very low mean firing rates but strong effect sizes. All methods were trained on a list of squared distances to a grid of place fields and on the rat head orientation, as described in methods. Far more even than the neocortical data, neurons were described much better by XGBoost and the ensemble method than by the GLM (**Figure 7B**). Many neurons shifted from being completely unpredictable by the GLM (pseudo- R^2 near zero) to very predictable by XGBoost and the ensemble (pseudo- R^2 above 0.2). These neurons thus have responses that do not correlate with firing in any one Gaussian place field. We note that the neural network performed poorly, likely due to the very low firing rates of most hippocampal cells (Supplementary Figure 2). The median spike rate of the 58 neurons in the dataset was just 0.2 spikes/s, and it was only on the four neurons with rates above 1 spikes/s that the neural network achieved pseudo- R^2 scores comparable to the GLM. The relative success of XGBoost was interesting given the failure of the neural network, and supported the general observation that boosted trees can work well with smaller and sparser datasets than those that neural networks generally require (Supplementary Figure 3). Thus for hippocampal cells, a method leveraging decision trees such as XGBoost or the ensemble is able to capture more structure in the neural response and thus demonstrate a deficiency of the parameterization of the GLM.

DISCUSSION

We analyzed the ability of various machine learning techniques at the task of predicting binned spike counts in three brain regions. We found that of the tested ML methods, XGBoost and the ensemble routinely predicted spike counts more accurately than did the GLM, which is a popular method for neural data. Feedforward neural networks did not always outperform the GLM and were often worse than XGBoost and the ensemble. Machine learning methods, especially LSTMs, also outperformed GLMs when covariate and spike history were included as inputs. The ML methods performed comparably well with and without feature engineering, even for the very low spike rates of the hippocampus dataset. These findings indicate that a standard ML approach can serve as a reliable benchmark to test if data meets the assumptions of a GLM. Furthermore, it may be quite common that standard ML outperforms GLMs given standard feature choices.

When a GLM fails to explain data as well as more expressive, nonlinear methods, the current parameterization of inputs must relate to the data with a different nonlinearity than is assumed by the GLM. Such situations have been identified several times in the literature (Butts et al., 2011; Freeman et al., 2015; Heitman et al., 2016; McIntosh et al., 2016). This unaccounted nonlinearity may produce feature weights that do not reflect true feature importance. A GLM will incorrectly predict no dependence on feature x whatsoever, for example, in the extreme case when the neural response to some feature x does not correlate with $\exp(x)$. The only way to ensure that feature weights can be reliably



interpreted is to find an input parameterization that maximizes the GLM's predictive power. ML methods can assist this process by indicating how much nonlinearity remains to be explained. New features can then be tested, such as those suggested by a search for maximally informative dimensions (Sharpee et al., 2004). In our analysis, then, the GLM underperforms because we have selected the suboptimal input features. It is always theoretically possible to linearize features such that a GLM obtains equal predictive power. ML methods can highlight the deficiency of features that might have otherwise seemed uncontroversial. When applying a GLM or any simple model to neural data, it is important to compare its predictive power with standard ML methods to ensure the neural response is properly understood.

There are other ways of estimating the performance of a method besides benchmark nonlinear methods. For example, if the same exact stimulus can be given many times in a row, then we can estimate neural variability without having to model how activity depends on stimulus features (Schoppe et al., 2016). This approach, however, requires that we can model how neural responses vary with repetition (Grill-Spector et al., 2006). This approach also makes it difficult to include spike history as an input, since the exact history is rarely repeated. We note that in some cases it may also be impossible to show the same stimulus multiple times, e.g., because eyes move. However, comparing these two classes of benchmark would be interesting on applications where both are feasible.

Advanced ML methods are not widely considered to be interpretable. Interpretation is not necessary for performance benchmarks, but it would be desirable to use these methods as standalone encoding models. We can better discuss this issue with a more precise definition of interpretability. Following Lipton, we make the distinction between a method's *post-hoc*

interpretability, the ease of justifying its predictions, and *transparency*, the degree to which its operation and internal parameters are human-readable or easily understandable (Lipton et al., 2016). A GLM is certainly more transparent than many ML methods due to its algorithmic simplicity. Certain nonlinear extensions of the GLM have also been designed to remain transparent (McFarland et al., 2013; Theis et al., 2013; Latimer et al., 2014; Williamson et al., 2015; Maheswaranathan et al., 2017). For high-level areas, though, such as V4, the linearized features may be difficult to be interpreted themselves (Yamins et al., 2014), though it may be possible to increase the interpretability of features (Kaardal et al., 2013). A GLM is also generally more conducive to *post-hoc* interpretations, though this is also possible with modern ML methods. It is possible, for example, to visualize the aspects of stimuli that most elicit a predicted response, as has been implemented in previous applications of neural networks to spike prediction (Lau et al., 2002; Prenger et al., 2004). Various other methods exist in the literature to enable *post-hoc* explanations (McAuley and Leskovec, 2013; Simonyan et al., 2013). Here we highlight Local Interpretable Model-Agnostic Explanations (LIME), an approach that fits simple models in the vicinity of single examples to allow a local interpretation (Ribeiro et al., 2016). On problems where interpretability is important, such capabilities for *post-hoc* justifications may prove sufficient.

Not all types of interpretability are necessary for a given task, and many scientific questions can be answered based on predictive ability alone. Questions of the form, "does feature x contribute to neural activity?" for example, or "is past activity necessary to explain current activity?" require no method transparency. One can simply ask whether predictive power increases with feature x 's inclusion or decreases upon its exclusion. Importance measures based on inclusion and

exclusion, or upon the strategy of shuffling a covariate of interest, are well-studied in statistics and machine learning (Bell and Wang, 2000; Strobl et al., 2008). Depending on the application, it may thus be worthwhile to ask not just whether different features could improve a GLM but also whether it is enough to use ML methods directly. It is possible for many questions to stay agnostic to the form of linearized features and directly use changes in predictive ability.

With ongoing progress in machine learning, many standard techniques are easy to implement and can even be automated. Ensemble methods, for example, remove the need to choose any one algorithm. Moreover, the choice of model-specific parameters is made easy by hyperparameter search methods and optimizers. We hope that this ease of use might encourage use in the neurosciences, thereby increasing the power and efficiency of studies involving neural prediction without requiring complicated, application-specific methods development (e.g., Corbett et al., 2012). Community-supported projects in automated machine learning, such as autoSklarn and auto-Weka, are quickly improving and promise to handle the entire regression workflow (Feurer et al., 2015; Kotthoff et al., 2016). Applied to neuroscience, these tools will allow researchers to gain descriptive power over current methods even with simple, out-of-the-box implementations.

Machine learning methods perform quite well and make minimal assumptions about the form of neural encoding. Models that seek to understand the form of the neural code can test if they systematically misconstrue the relationship between stimulus and response by comparing their performance to these benchmarks. Encoding models built with machine learning can thus greatly aid the construction of models that capture arbitrary nonlinearity and more accurately describe neural activity.

REFERENCES

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., et al. (2016). Tensorflow: large-scale machine learning on heterogeneous distributed systems. *arXiv:1603.04467* [preprint].
- Aljadeff, J., Lansdell, B. J., Fairhall, A. L., and Kleinfeld, D. (2016). Analysis of neuronal spike trains, deconstructed. *Neuron* 91, 221–259. doi: 10.1016/j.neuron.2016.05.039
- Amirikian, B., and Georgopoulos, A. P. (2000). Directional tuning profiles of motor cortical cells. *Neurosci. Res.* 36, 73–79. doi: 10.1016/S0168-0102(99)00112-1
- BayesianOptimization (2016). *GitHub Repository*.
- Bell, D. A., and Wang, H. (2000). A formalism for relevance and its application in feature subset selection. *Mach. Learn.* 41, 175–195. doi: 10.1023/A:1007612503587
- Brown, E. N., Frank, L. M., Tang, D., Quirk, M. C., and Wilson, M. A. (1998). A statistical paradigm for neural spike train decoding applied to position prediction from ensemble firing patterns of rat hippocampal place cells. *J. Neurosci.* 18, 7411–7425. doi: 10.1523/JNEUROSCI.18-18-07411.1998
- Butts, D. A., Weng, C., Jin, J., Alonso, J. M., and Paninski, L. (2011). Temporal precision in the visual pathway through the interplay of excitation and stimulus-driven suppression. *J. Neurosci.* 31, 11313–11327. doi: 10.1523/JNEUROSCI.0434-11.2011
- Cameron, A. C., and Windmeijer, F. A. (1997). An R-squared measure of goodness of fit for some common nonlinear regression models. *J. Econom.* 77, 329–342.
- Chen, T., and Guestrin, C. (2016). Xgboost: a scalable tree boosting system. *arXiv preprint arXiv:1603.02754* [preprint].
- Chichilnisky, E. (2001). A simple white noise analysis of neuronal light responses. *Netw. Comp. Neural Syst.* 12, 199–213. doi: 10.1080/713663221
- Chollet, F. (2015). Keras. *GitHub repository*. Available Online at: <https://github.com/keras-team/keras>
- Corbett, E. A., Perreault, E. J., and Kording, K. P. (2012). Decoding with limited neural data: a mixture of time-warped trajectory models for directional reaches. *J. Neural Eng.* 9:036002. doi: 10.1088/1741-2560/9/3/036002
- Domencich, T. A., and McFadden, D. (1975). *Urban Travel Demand-A Behavioral Analysis*, Oxford: North-Holland Publishing Company Limited.
- Fernandes, H. L., Stevenson, I. H., Phillips, A. N., Segraves, M. A., and Kording, K. P. (2014). Saliency and saccade encoding in the frontal eye field during natural scene search. *Cereb. Cortex* 24, 3232–3245. doi: 10.1093/cercor/bht179
- Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., and Hutter, F. (2015). “Efficient and robust automated machine learning,” in *Advances in Neural Information Processing Systems*, Vol. 28, eds C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (Curran Associates, Inc.), 2962–2970. Available online at: <http://papers.nips.cc/paper/5872-efficient-and-robust-automated-machine-learning.pdf>
- Freeman, J., Field, G. D., Li, P. H., Greschner, M., Gunning, D. E., Mathieson, K., et al. (2015). Mapping nonlinear receptive field structure in primate retina at single cone resolution. *Elife* 4:e05241. doi: 10.7554/eLife.05241
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Ann. Statist.* 29, 1189–1232. doi: 10.1214/aos/1013203451
- Friedman, J., Hastie, T., and Tibshirani, R. (2000). Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *Ann. Statist.* 28, 337–407. doi: 10.1214/aos/1016120463

The code used for this publication is available at <https://github.com/KordingLab/spykesML>. We invite researchers to adapt it freely for future problems of neural prediction.

AUTHOR CONTRIBUTIONS

KK and HF first conceived the project. TT, CV, and RC gathered and curated macaque data. AB prepared the manuscript and performed the analyses, for which HF and PR assisted. LM and KK supervised, and all authors assisted in editing.

FUNDING

LM acknowledges the following grants from the National Institute of Neurological Disorders and Stroke (<https://www.ninds.nih.gov/>): NS074044, NS048845, NS053603, NS095251. CV recognizes support from the Biomedical Data Driven Discovery (BD3) Training Program funded through the National Institute of Health, grant number 5T32LM012203-02. KK acknowledges support from National Institute of Health (<https://www.nih.gov/>) grants R01NS063399, R01NS074044, and MH103910. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. RC thanks NSF GRFP DGE-1324585.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fncom.2018.00056/full#supplementary-material>

- Gers, F. A., Schmidhuber, J., and Cummins, F. (2000). Learning to forget: Continual prediction with LSTM. *Neural Comput.* 12, 2451–2471. doi: 10.1162/089976600300015015
- Gerwinn, S., Macke, J. H., and Bethge, M. (2010). Bayesian inference for generalized linear models for spiking neurons. *Front. Comp. Neurosci.* 4:12. doi: 10.3389/fncom.2010.00012
- Grill-Spector, K., Henson, R., and Martin, A. (2006). Repetition and the brain: neural models of stimulus-specific effects. *Trends Cogn. Sci.* 10, 14–23. doi: 10.1016/j.tics.2005.11.006
- Heitman, A., Brackbill, N., Greschner, M., Sher, A., Litke, A. M., and Chichilnisky, E. (2016). Testing pseudo-linear models of responses to natural scenes in primate retina. *bioRxiv:045336* [preprint]. doi: 10.1101/045336
- Ho, T. K. (1998). The random subspace method for constructing decision forests. *IEEE Trans. Pattern Anal. Mach. Intell.* 20, 832–844.
- Kaardal, J., Fitzgerald, J. D., Berry, M. J., and Sharpee, T. O. (2013). Identifying functional bases for multidimensional neural computations. *Neural Comput.* 25, 1870–1890. doi: 10.1162/NECO_a_00465
- Kaggle Winner's Blog (2016). Available Online at: <http://blog.kaggle.com/>
- Kotthoff, L., Thornton, C., Hoos, H. H., Hutter, F., and Leyton-Brown, K. (2016). Auto-WEKA 2.0: automatic model selection and hyperparameter optimization in WEKA. *J. Mach. Learn. Res.* 17, 1–5. doi: 10.1145/2487575.2487629
- Latimer, K. W., Chichilnisky, E., Rieke, F., and Pillow, J. W. (eds). (2014). “Inferring synaptic conductances from spike trains with a biophysically inspired point process model,” in *Advances in Neural Information Processing Systems* (Curran Associates, Inc.), 954–962.
- Lau, B., Stanley, G. B., and Dan, Y. (2002). Computational subunits of visual cortical neurons revealed by artificial neural networks. *Proc. Natl. Acad. Sci. U.S.A.* 99, 8974–8979. doi: 10.1073/pnas.122173799
- Lipton, Z. C. (2016). The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*.
- Maheswaranathan, N., Baccus, S. A., and Ganguli, S. (2017). Inferring hidden structure in multilayered neural circuits. *bioRxiv: 120956* [preprint]. doi: 10.1101/120956
- McAuley, J., and Leskovec, J. (eds). (2013). “Hidden factors and hidden topics: understanding rating dimensions with review text,” in *Proceedings of the 7th ACM Conference on Recommender Systems* (Hong Kong: ACM).
- McFarland, J. M., Cui, Y., and Butts, D. A. (2013). Inferring nonlinear neuronal computation based on physiologically plausible inputs. *PLoS Comput. Biol.* 9:e1003143. doi: 10.1371/journal.pcbi.1003143
- McIntosh, L., Maheswaranathan, N., Nayebi, A., Ganguli, S., Baccus, S. (2016). “Deep learning models of the retinal response to natural scenes,” in *Advances in Neural Information Processing Systems, Vol. 29*, eds D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett (Curran Associates, Inc.), 1369–1377. Available online at: <http://papers.nips.cc/paper/6388-deep-learning-models-of-the-retinal-response-to-natural-scenes.pdf>
- Mizuseki, K., Sirota, A., Pastalkova, E., and Buzsáki, G. (2009a). Multi-unit recordings from the rat hippocampus made during open field foraging. doi: 10.6080/K0Z60KZ9
- Mizuseki, K., Sirota, A., Pastalkova, E., and Buzsáki, G. (2009b). Theta oscillations provide temporal windows for local circuit computation in the entorhinal-hippocampal loop. *Neuron* 64, 267–280. doi: 10.1016/j.neuron.2009.08.037
- Nelder, J. A., and Baker, R. J. (1972). Generalized linear models. *Encyclop. Statist. Sci.* 135, 370–384. doi: 10.2307/2344614
- Paninski, L. (2004). Maximum likelihood estimation of cascade point-process neural encoding models. *Netw. Comp. Neural Sys.* 15, 243–262. doi: 10.1088/0954-898X_15_4_002
- Paninski, L., Fellows, M. R., Hatsopoulos, N. G., and Donoghue, J. P. (2004a). Spatiotemporal tuning of motor cortical neurons for hand position and velocity. *J. Neurophysiol.* 91, 515–532. doi: 10.1152/jn.00587.2002
- Paninski, L., Shoham, S., Fellows, M. R., Hatsopoulos, N. G., and Donoghue, J. P. (2004b). Superlinear population encoding of dynamic hand trajectory in primary motor cortex. *J. Neurosci.* 24, 8551–8561. doi: 10.1523/JNEUROSCI.0919-04.2004
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.* 12, 2825–2830.
- Pillow, J. W., Paninski, L., Uzzell, V. J., Simoncelli, E. P., and Chichilnisky, E. (2005). Prediction and decoding of retinal ganglion cell responses with a probabilistic spiking model. *J. Neurosci.* 25, 11003–11013. doi: 10.1523/JNEUROSCI.3305-05.2005
- Pillow, J. W., Shlens, J., Paninski, L., Sher, A., Litke, A. M., Chichilnisky, E., et al. (2008). Spatio-temporal correlations and visual signalling in a complete neuronal population. *Nature* 454, 995–999. doi: 10.1038/nature07140
- Prenger, R., Wu, M. C., David, S. V., and Gallant, J. L. (2004). Nonlinear V1 responses to natural scenes revealed by neural network analysis. *Neural Netw.* 17, 663–679. doi: 10.1016/j.neunet.2004.03.008
- Prud'homme, M. J., and Kalaska, J. F. (1994). Proprioceptive activity in primate primary somatosensory cortex during active arm reaching movements. *J. Neurophysiol.* 72, 2280–2301. doi: 10.1152/jn.1994.72.5.2280
- Ramkumar, P., Jas, M., Achakulvisut, T., Idrizović, A., themantolope, Acuna, D. E., et al. (2017). *Pyglnnet 1.0.1*. (Chicago, IL).
- Ramkumar, P., Lawlor, P. N., Glaser, J. L., Wood, D. K., Phillips, A. N., Segraves, M. A., et al. (2016). Feature-based attention and spatial selection in frontal eye fields during natural scene search. *J. Neurophysiol.* 116, 1328–1343. doi: 10.1152/jn.01044.2015
- Ribeiro, M. T., Singh, S., and Guestrin, C. (eds). (2016). “Why should i trust you?: Explaining the predictions of any classifier,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Francisco, CA: ACM).
- Saleh, M., Takahashi, K., and Hatsopoulos, N. G. (2012). Encoding of coordinated reach and grasp trajectories in primary motor cortex. *J. Neurosci.* 32, 1220–1232. doi: 10.1523/JNEUROSCI.2438-11.2012
- Schmidhuber, J. (2015). Deep learning in neural networks: an overview. *Neural Netw.* 61:85–117. doi: 10.1016/j.neunet.2014.09.003
- Schoppe, O., Harper, N. S., Willmore, B. D., King, A. J., and Schnupp, J. W. (2016). Measuring the performance of neural models. *Front. Comput. Neurosci.* 10:10. doi: 10.3389/fncom.2016.00010
- Schwartz, O., Pillow, J. W., Rust, N. C., and Simoncelli, E. P. (2006). Spike-triggered neural characterization. *J. Vis.* 6:13. doi: 10.1167/6.4.13
- Sharpee, T., Rust, N. C., and Bialek, W. (2004). Analyzing neural responses to natural signals: maximally informative dimensions. *Neural Comput.* 16, 223–250. doi: 10.1162/089976604322742010
- Simoncelli, E. P., Paninski, L., Pillow, J., and Schwartz, O. (2004). “Characterization of neural responses with stochastic stimuli,” in *The cognitive neurosciences, 3rd edn* ed M. Gazzaniga (MIT Press), 327–338.
- Simonyan, K., Vedaldi, A., and Zisserman, A. (2013). Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv:1312.6034*[preprint].
- Snoek, J., Larochelle, H., and Adams, R. P. (2012). “Practical Bayesian optimization of machine learning algorithms,” in *Advances in Neural Information Processing Systems, Vol. 25*, eds F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (Curran Associates, Inc.), 2951–2959. Available online at: <http://papers.nips.cc/paper/4522-practical-bayesian-optimization-of-machine-learning-algorithms.pdf>
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* 15, 1929–1958.
- Stevenson, I. H., Cherian, A., London, B. M., Sachs, N. A., Lindberg, E., Reimer, J., et al. (2011). Statistical assessment of the stability of neural movement representations. *J. Neurophysiol.* 106, 764–774. doi: 10.1152/jn.00626.2010
- Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., and Zeileis, A. (2008). Conditional variable importance for random forests. *BMC Bioinformatics* 9:307. doi: 10.1186/1471-2105-9-307
- Team, T. T. D., Al-Rfou, R., Alain, G., Almahairi, A., Angermueller, C., Bahdanau, D., et al. (2016). Theano: a Python framework for fast computation of mathematical expressions. *arXiv: 160502688*[preprint].

- Theis, L., Chagas, A. M., Arnstein, D., Schwarz, C., and Bethge, M. (2013). Beyond GLMs: a generative mixture modeling approach to neural system identification. *PLoS Comput. Biol.* 9:e1003356. doi: 10.1371/journal.pcbi.1003356
- Truccolo, W., Eden, U. T., Fellows, M. R., Donoghue, J. P., and Brown, E. N. (2005). A point process framework for relating neural spiking activity to spiking history, neural ensemble, and extrinsic covariate effects. *J. Neurophysiol.* 93, 1074–1089. doi: 10.1152/jn.00697.2004
- Weber, A. I., and Pillow, J. W. (2016). Capturing the dynamical repertoire of single neurons with generalized linear models. *arXiv: 160207389*[preprint].
- Williamson, R. S., Sahani, M., and Pillow, J. W. (2015). The equivalence of information-theoretic and likelihood-based methods for neural dimensionality reduction. *PLoS Comput. Biol.* 11:e1004141. doi: 10.1371/journal.pcbi.1004141
- Wolpert, D. H. (1992). Stacked generalization. *Neural Netw.* 5, 241–259.
- Wu, M. C. K., David, S. V., and Gallant, J. L. (2006). Complete functional characterization of sensory neurons by system identification. *Annu. Rev. Neurosci.* 29, 477–505. doi: 10.1146/annurev.neuro.29.051605.113024
- Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., and DiCarlo, J. J. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceed. Natl. Acad. Sci. U.S.A.* 111, 8619–8624. doi: 10.1073/pnas.1403112111

Conflict of Interest Statement: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2018 Benjamin, Fernandes, Tomlinson, Ramkumar, VerSteeg, Chowdhury, Miller and Kording. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.