

Learning low-dimensional inputs for brain-machine interface control

Jorge A. Menéndez^{1,2} & Peter E. Latham¹

¹ Gatsby Computational Neuroscience Unit, ² CoMPLEX; University College London

A critical, but poorly understood, feature of the motor system is that it can learn new motor behaviors without forgetting old ones. A remarkable demonstration of this ability is given by brain-machine interfaces (BMI), which show that primates are capable of learning experimenter-imposed mappings from motor cortical activity to movement. Previous models of this phenomenon (Legenstein et al., 2010; Waernberg & Kumar, 2017) invoke synaptic plasticity in the local network. However, several observations are at odds with this hypothesis. In some cases, the correlation structure in the local population is largely conserved between BMI and manual control (Hwang et al., 2013) and between different BMI decoders (Golub et al., 2018), suggesting that the local circuitry has not undergone any changes. Furthermore, primates are able to remember and rapidly switch between previously used decoders (Ganguly & Carmena, 2009) – something that local plasticity rules have considerable difficulty achieving. Motivated by these observations, here we propose an alternative hypothesis: learning is isolated to upstream commands driving the local motor cortical circuit. Under this hypothesis, efficient learning can be achieved by restricting the space of possible upstream commands to those used during natural movement, which has been previously suggested to underlie BMI control. We formalize this hypothesis as a control problem in which upstream commands drive a recurrent network model of motor cortex. The formalism provides a quantitative relationship between the experimenter-determined BMI decoder and the difficulty of learning the task; we show that this relationship can account for previous observations on learning different BMI decoders (Sadler et al., 2014). Moreover, we argue that solving the control problem in this way makes efficient learning viable without computing gradients - a critical feature of BMI learning.

Additional Detail: We consider a simple rate-based model of the local motor cortical population $\mathbf{x} \in \mathbb{R}^n$, which receives input from an upstream population $\mathbf{u} \in \mathbb{R}^m$,

$$\dot{\mathbf{x}} = -\mathbf{x} + \mathbf{W}\phi(\mathbf{x}) + \mathbf{B}\mathbf{u}.$$

We assume that the recurrent and feedforward weights ensure that when the network is driven by an input $\mathbf{u} = f(\mathbf{v})$, the resulting motor cortical activity at time t , $\mathbf{x}(t; \mathbf{v})$, leads to a reaching movement in the direction \mathbf{v} .

Under BMI control, however, the same pattern of activity may lead to a reach in a completely different direction $\mathbf{v}^{\text{BMI}}(t) = \mathbf{D}\mathbf{x}(t; \mathbf{v})$, given an arbitrary linear decoder \mathbf{D} (as is used in practice). To compensate for this new mapping from neural activity to reaching direction, one could modify the network in several ways: (1) modify the recurrent synaptic circuitry by changing the n^2 synaptic weights W_{ij} ; (2) modify the synapses from the upstream population by changing the nm synaptic weights B_{ij} ; or (3) modify the m upstream input activities u_i . With no gradient information available, learning in a high dimensional space is slow (Werfel, Xie & Seung, 2004 NIPS), so the natural approach is to learn the m inputs u_i . However, even m is typically large.

An alternative approach would be to exploit the circuitry already present for reaching, and instead fix $\mathbf{u} = f(\tilde{\mathbf{v}})$ and learn the appropriate directional signal $\tilde{\mathbf{v}}$ that will drive the network to the desired state,

$$\tilde{\mathbf{v}} = \arg \min_{\mathbf{v}} E(\mathbf{D}\mathbf{x}(t; \mathbf{v}), \mathbf{v}^*)$$

where E measures the error between the desired reach direction, \mathbf{v}^* , and the actual reach direction $\mathbf{v}^{\text{BMI}}(t) = \mathbf{D}\mathbf{x}(t; \mathbf{v})$. This corresponds to a “re-aiming” strategy in which the primate aims in the direction $\tilde{\mathbf{v}}$ to achieve a reach in the direction \mathbf{v}^* . This reduces the learning problem to two dimensions, allowing rapid learning without explicit gradients. At the same time, however, restricting the inputs \mathbf{u} to live in a low-dimensional space might prevent us from finding a biophysically feasible solution to the control problem. In particular, the input required to attain a particular reach trajectory, $\mathbf{u} = f(\tilde{\mathbf{v}})$, could require unrealistically high firing rates.

So when does a suitable solution exist? We answer this question in the context of the simple control problem of ensuring that the BMI readout at some future time $\mathbf{D}\mathbf{x}_f(\mathbf{v}) \equiv \mathbf{D}\mathbf{x}(t_f; \mathbf{v})$ take on a desired state \mathbf{v}^* ,

$$E(\mathbf{D}\mathbf{x}_f(\mathbf{v}), \mathbf{v}^*) = \|\mathbf{D}\mathbf{x}_f(\mathbf{v}) - \mathbf{v}^*\|^2 + \frac{\gamma}{m} \|\mathbf{u}\|^2$$

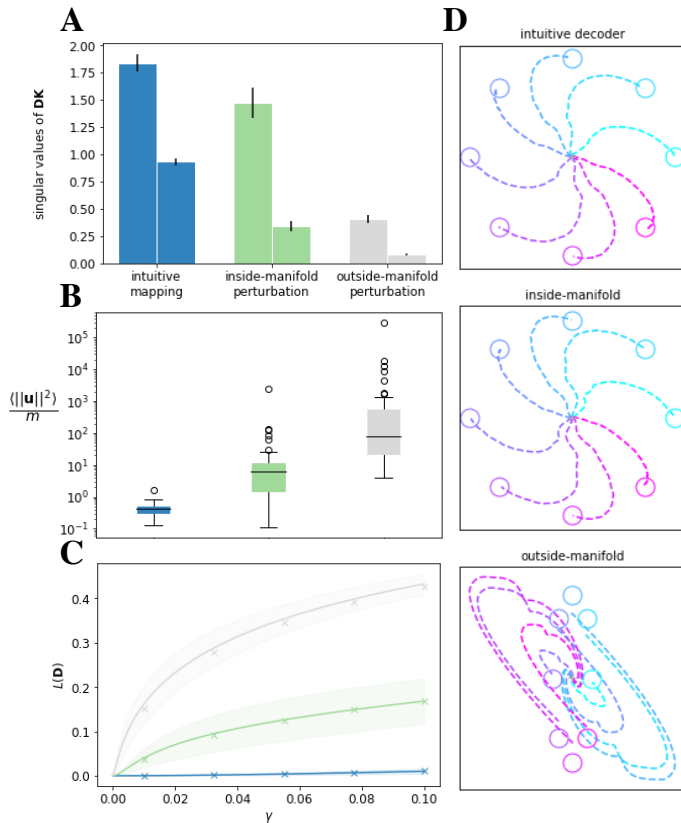
where γ specifies a cost on the mean squared input activity, which must fall in a physiologically reasonable range. Inspired by observed cosine tuning to reach direction in motor cortex (Georgopoulos et al., 1982), we let the

inputs depend linearly on the reach direction by setting $f(\mathbf{v}) = \mathbf{M}\mathbf{v}$, where the elements of \mathbf{M} are taken to be random and *iid*. In the case of linear dynamics, $\phi(\mathbf{x}) = \mathbf{x}$, we can solve for $\mathbf{x}(t)$ to obtain $\mathbf{x}_f(\mathbf{v}) = \mathbf{K}\mathbf{v}$, where \mathbf{K} depends on the parameters of the network \mathbf{W} , \mathbf{B} and \mathbf{M} . This results in a unique solution to the above optimization problem, of the form $\tilde{\mathbf{v}} = \mathbf{G}\mathbf{v}^*$ where \mathbf{G} depends on \mathbf{D} , \mathbf{K} and γ . Crucially, we can use this expression to evaluate a given decoder \mathbf{D} by quantifying the average error achieved when using this optimal input,

$$L(\mathbf{D}) = \langle \|\mathbf{D}\mathbf{x}_f(\tilde{\mathbf{v}}) - \mathbf{v}^*\|^2 \rangle_{\mathbf{v}^*} = \langle \|\mathbf{D}\mathbf{K}\mathbf{G}\mathbf{v}^* - \mathbf{v}^*\|^2 \rangle_{\mathbf{v}^*} = \frac{\gamma^2}{2} \sum_{i=1}^2 \frac{1}{(s_i^2 + \gamma)^2}$$

where the average is over intended reach directions \mathbf{v}^* uniformly distributed on the unit circle and the s_i are the singular values of the matrix \mathbf{DK} .

To achieve low error with γ sufficiently large to ensure that \mathbf{u} is in a biophysical range, the singular values, s_i , must be large. Crucially, we note that the elements of the matrix \mathbf{DK} are dot products in n -dimensional space: $(\mathbf{DK})_{ij} = \sum_{k=1}^n D_{ik}K_{kj}$. The magnitude of the singular values will thus depend on whether these dot products are large, which in turn depends on the alignment between the rows and columns of \mathbf{D} and \mathbf{K} , respectively. More precisely, if the rows of \mathbf{D} and columns of \mathbf{K} are aligned, their dot products will be $\mathcal{O}(\sqrt{n})$ larger than if they are random. Given that the number of neurons in the network n is large, this will constitute a significant difference in the singular values s_i . The matrix \mathbf{K} is intimately related to the network dynamics, so ultimately this analysis reveals that the decoder must be aligned with the natural dynamics of the motor cortical circuit for a solution to the above control problem to be possible with reasonably strong inputs.



A. Singular values of matrix \mathbf{DK} for each decoder. **B.** Magnitude of optimal inputs $\mathbf{u} = \mathbf{M}\mathbf{G}\mathbf{v}^*$. **C.** Decoder loss evaluated for different input costs γ . Lines denote the theoretical values and x's denote empirical measurements from simulations. **D.** Example reaches from one simulated experiment, using optimal inputs for $\gamma = .001$. In A,B,C, error bars denote standard error of the mean over 50 random network initializations, with $W_{ij} \sim \mathcal{N}(0, 1/n)$, $B_{ij} \sim \mathcal{N}(0, 1/m)$, $M_{ij} \sim \mathcal{N}(0, 1/2)$, $n = 2000$, $m = 500$. Decoders were fit following the procedure of Sadtler et al., using neural activity simulated with the true target directions as input, $\mathbf{u} = \mathbf{M}\mathbf{v}^*$. We fit the decoders to a random subset of 100 neurons from the full network.

An elegant study by Sadtler et al. (2014, Nature) asked whether decoders were easier to learn if they were aligned with the local network dynamics. To test whether our theory could account for their results, we evaluated each of the decoders they used using the above defined decoder loss $L(\mathbf{D})$. We show that indeed only those decoders that are aligned with the network dynamics (the “intuitive mapping” and “inside-manifold perturbation”) allow for solutions with reasonably strong inputs (fig. B), and that this depends on the singular values of \mathbf{DK} exactly as predicted by the theory (fig. A,C). Indeed, Sadtler et al. found that primates were able to rapidly learn these decoders, whereas the other (“outside-manifold perturbation”) required more time to learn. Our theory offers an interpretation of these results as a consequence of learning in different spaces: whereas the “intuitive” and “inside-manifold” decoders can be learnt through optimization in the two-dimensional space occupied by the inputs used during manual reaching, learning to use the “outside-manifold” decoder requires optimizing over higher dimensions. Given that no gradient is available for this optimization, learning will be substantially slower in this case.

Even in the “inside-manifold” case, however, certain decoders may be easier to learn than others. The above framework provides a way to evaluate how this might depend on the statistics of the decoder \mathbf{D} and of the tuning curves imposed by \mathbf{M} . Alternatively, this might depend on the nature of the optimal visuomotor transformation \mathbf{G} . These directions are being pursued.