# Statistical physics, Bayesian inference and neural information processing

Erin Grant

*Gatsby Computational Neuroscience Unit & Sainsbury Wellcome Centre, UCL, London, UK*

Sandra B. Nestler

*Institute of Neuroscience and Medicine (INM-6) and Institute for Advanced Simulation (IAS-6) and*
*JARA-Institute Brain Structure-Function Relationships (INM-10), Jülich Research Centre, Jülich, Germany and*
*Department of Physics, Faculty 1, RWTH Aachen University, Aachen, Germany*

Berfin Şimşek

*Chair of Statistical Field Theory (CSFT) & Laboratory of Computational Neuroscience (LCN), École Polytechnique Fédérale de Lausanne*

Sara A. Solla

*Department of Neuroscience, Northwestern University, Chicago, IL 60611, USA and*
*Department of Physics and Astronomy, Northwestern University, Evanston, IL 60208, USA*

**CONTENTS**

## I. LECTURE 1: STATISTICAL PHYSICS, BAYESIAN INFERENCE, AND NEURAL INFORMATION PROCESSING

The brain interprets the world; an embodied brain causes actions that change the world. The brain interacts with an external dynamical system, the world, that follows causal relationships, as in $\vec{y}$ follows $\vec{x}$. However, the brain

$$\vec{x} \longrightarrow \vec{y}$$

receives the same input $\vec{x}$, processes it, and acts in a manner that affects the output. This interactive process be-
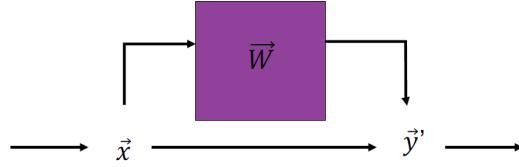


FIG. 1. The brain as part of the world: receiving inputs from the world and producing outputs that change the state of the world.

tween the brain and the world can be framed in terms of input-output maps between high-dimensional state vectors $\vec{x} = \{x_1, x_2, \ldots, x_N\} \to \vec{y} = \{y_1, y_2, \ldots, y_R\}$. The output of the system is a function of its input,

$$\vec{y} = f(\vec{x}) . \tag{1}$$

### A. Learning to model an input-output map

We are particularly interested in the case in which the input-output map can be modeled as implemeted by a network of *neurons* characterized by connectivity parameters $\vec{W}$, such that

$$\vec{y} = f_{\vec{W}}(\vec{x}) . \tag{2}$$



FIG. 2. The model as an input-output-map specified by its parameters $W$.

What specifies the value of the parameters $\vec{W}$? The parameters of this network are determined by a learning mechanism that relies on data $\vec{\xi}$ that provides information about the desired input-output map:

$$\vec{\xi}^{\mu} = (\vec{x}^{\mu}, \vec{y}^{\mu}) \quad 1 \leq \mu \leq m . \tag{3}$$

These $m$ input-output pairs may be exact or corrupted by noise, but they are examples of the desired map. Supervised learning can be applied to infer a map that approximates the desired solution.

The comparison between the network's output and the desired output defines an error, and the parameters are adapted to reduce this error. Given an example of the desired map, the error made by a specific module $\vec{W}$ on this example is

$$E(\vec{W} \mid \vec{x}, \vec{y}) = d\left(\vec{y}, f_{\vec{W}}(\vec{x})\right) = \frac{1}{2}\left(\vec{y} - f_{\vec{W}}(\vec{x})\right)^2 . \tag{4}$$

The loss function for such training does not need to be quadratic; it is chosen to be some kind of distance defined in output space. The choice of error function does not need to satisfy all requirements that define a proper distance: it does not need
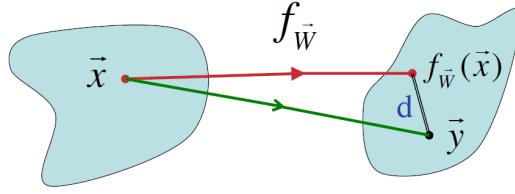
FIG. 3. To evaluate a model, compare the output $f_{\vec{W}}(\vec{x})$ predicted by the model to the desired output $\vec{y}$.

to be symmetric and it does not need to obey the triangle inequality. It only needs to satisfy that it is equal to zero if and only if its two arguments are equal.

Given a training set of size $m$

$$\vec{\xi}^{\mu} = (\vec{x}^{\mu}, \vec{y}^{\mu}) \quad 1 \leq \mu \leq m \,, \tag{5}$$

we construct a cost function that measures the average error over the training set, the *learning error*:

$$E_L(\vec{W}) = (1/m) \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{x}^{\mu}, \vec{y}^{\mu}\right) \,. \tag{6}$$

Most learning algorithms are based on finding the parameters $\vec{W}^*$ that minimize this learning error. Most often, this is achieved by a form of gradient descent, where the weights are modified in small steps that depend in direction and size on the gradient of the learning error. The system will then eventually converge to a local minimum of the loss landscape.



FIG. 4. Loss landscape with a global and local minima.

## B.  Perceptron learning by gradient descent

A perceptron is a network model which obtains its scalar output from an $N$-dimensional input through a linear mapping followed by a soft nonlinearity:

$$y = g\left(\sum_{i=1}^{N} w_i x_i + w_0\right) = g\left(\vec{w}^T \vec{x}\right) \,, \tag{7}$$

where we have extended the input to $N + 1$ dimensions by adding a component $x_0 = 1$ to account for the bias $w_0$. The error on the $\mu$-th example and its gradient are then given by

$$E^{\mu} = \frac{1}{2}\left(y^{\mu} - g\left(\sum_{i=1}^{N} w_i x_i^{\mu} + w_0\right)\right)^2 \,, \tag{8}$$

$$\frac{\partial E^{\mu}}{\partial w_i} = -\left(y^{\mu} - g\left(\vec{w}^T \vec{x}^{\mu}\right)\right) g'\left(\vec{w}^T \vec{x}^{\mu}\right) x_i^{\mu} \,. \tag{9}$$

We can then formulate the gradient descent learning as a delta rule [1]. To this end, we define the shorthand $\delta^\mu$ as in

$$\frac{\partial E^\mu}{\partial w_i} = -\left(y^\mu - g\left(\vec{w}^T \vec{x}^\mu\right)\right) g'\left(\vec{w}^T \vec{x}^\mu\right) x_i^\mu \equiv -\delta^\mu x_i^\mu$$

$$\delta^\mu \equiv g'\left(\vec{w}^T \vec{x}^\mu\right)\left(y^\mu - g\left(\vec{w}^T \vec{x}^\mu\right)\right) \tag{10}$$

$$w_i \to w_i + \Delta w_i = w_i - \eta \frac{\partial E^\mu}{\partial w_i} = w_i + \eta \delta^\mu x_i^\mu \, .$$

The weight update for the perceptron is then given by

$$\Delta w_i^\mu = \eta \delta^\mu x_i^\mu \, , \tag{11}$$

where $\eta$ is a learning rate.

## C. Configuration space

Learning in artificial neural networks can also be regarded from an ensemble viewpoint [2, 3]. The ensemble of all possible networks compatible with a given architecture is described by its configuration space $\{\vec{W}\}$ . This space defines all possible functions that the network can implement given its architecture.



FIG. 5. A generic configuration space.

The weights $\vec{W}$ follow from a normalized prior distribution $\rho_0(\vec{W})$:

$$\int \rho_0(\vec{W}) d\vec{W} = 1. \tag{12}$$

This prior distribution assumes no information about the data, and should be chosen to be as unrestricted as possible.

### 1. Error-free learning

For each example $\vec{\xi}^\mu = (\vec{x}^\mu, \vec{y}^\mu)$ in the training set, define a masking function $\Theta$ as

$$\Theta\left(\vec{W}, \vec{\xi}^\mu\right) = \begin{cases} 1, & \text{if } f_{\vec{W}}\left(\vec{x}^\mu\right) = \vec{y}^\mu \\ 0, & \text{otherwise} \end{cases} \tag{13}$$

As we iteratively draw samples, the masking function $\Theta$ will eliminate all network configurations that do not satisfy the constraints imposed by the samples. The probability distribution is modified multiplicatively:

$$\rho_0(\vec{W})$$
$$\Rightarrow \rho_0(\vec{W}) \, \Theta(\vec{W}, \vec{\xi}^1) \tag{14}$$
$$\Rightarrow \rho_0(\vec{W}) \, \Theta(\vec{W}, \vec{\xi}^1) \, \Theta(\vec{W}, \vec{\xi}^2)$$

After the network has seen all the examples in the training set, the configuration space shrinks to

$$Z_m = \int d\vec{W} \rho_0(\vec{W}) \prod_{\mu=1}^{m} \Theta(\vec{W}, \vec{\xi}^\mu) \, , \tag{15}$$

FIG. 6. The configuration space shrinks during learning. Configurations that do not satisfy the constraints imposed by the examples are masked away.

with

$$Z_m \leq Z_{m-1} \leq \ldots \leq Z_1 \leq Z_0 = 1, \tag{16}$$

as with each iteration an additioanl subset of the configuration space may be eliminated.

### 2. Learning from noisy data

Consider an error on the $\mu$th example:
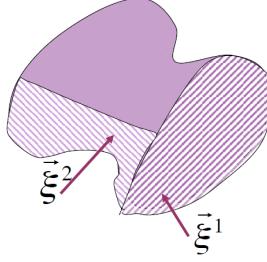
$$E\left(\vec{W} \mid \vec{\xi}^{\mu}\right) = d\left(\vec{y}^{\mu}, f_{\vec{W}}\left(\vec{x}^{\mu}\right)\right) \tag{17}$$

If $f_{\vec{W}}\left(\vec{x}^{\mu}\right) = \vec{y}^{\mu}$, the error function vanishes and the corresponding weights are not masked away: $E\left(W \mid \vec{\xi}^{\mu}\right) = 0 \Rightarrow \Theta\left(\vec{W}, \vec{\xi}^{\mu}\right) = 1$

If $f_{\vec{W}}\left(\vec{x}^{\mu}\right) \neq \vec{y}^{\mu}$, instead of setting $\Theta\left(\vec{W}, \vec{\xi}^{\mu}\right) = 0$ we can introduce a survival probability.

$$\Theta\left(\vec{W}, \vec{\xi}^{\mu}\right) \rightarrow \exp\left(-\beta E\left(\vec{W} \mid \vec{\xi}^{\mu}\right)\right) . \tag{18}$$

This corresponds to the difference between hard and soft masking, where configurations are attenuated by a factor exponentially controlled by the error made on the data instead of being eliminated. The exact choice of survival probability seems ad hoc at this point, but will become rigorously justified.



FIG. 7. Soft masking: Instead of eliminating subsets of configuration space, these regions are attenuated by an exponential factor.

The configuration space changes multiplicatively through a product of exponentials. The probability density becomes

$$\begin{aligned}
\rho_0(\vec{W}) &\Rightarrow \rho_0(\vec{W}) \exp\left(-\beta E\left(\vec{W} \mid \vec{\xi}^1\right)\right) \\
&\Rightarrow \rho_0(\vec{W}) \exp\left(-\beta E\left(\vec{W} \mid \vec{\xi}^1\right)\right) \exp\left(-\beta E\left(\vec{W} \mid \vec{\xi}^2\right)\right) \rightarrow \ldots
\end{aligned} \tag{19}$$

Given the full training set of size $m$, the effective volume of the configuration space becomes

$$Z_m = \int d\vec{W} \rho_0(\vec{W}) \prod_{\mu=1}^{m} \exp\left(-\beta E\left(\vec{W} \mid \vec{\xi}^\mu\right)\right) \tag{20}$$

Combining the product of exponentials into the exponential of the sum, we recover the mean error over all presented examples:

$$Z_m = \int d\vec{W} \rho_0(\vec{W}) \exp\left(-m\beta E_L(\vec{W})\right), \tag{21}$$

with

$$E_L(\vec{W}) = (1/m) \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{\xi}^\mu\right). \tag{22}$$

This formulation, based on the learning error and explicitly showing the size $m$ of the training set in the exponent, will allow us to establish a helpful connection to statistical physics.

### 3.  Thermodynamics of learning

In the ensemble description, the ensemble of all possible networks is described by the prior density $\rho_0(W)$, and the ensemble of trained networks is described by the posterior density $\rho_m(\vec{W})$ :

$$\rho_m(\vec{W}) = \frac{1}{Z_m} \rho_0(\vec{W}) \exp\left(-\beta m E_L(\vec{W})\right) \tag{23}$$

The equation for $\rho_m(\vec{W})$ takes the form of a Gibbs distribution. In comparison to statistical physics, the noise control parameter $\beta$ plays the role of the inverse temperature and the learning error plays the role of an energy function; the latter is extensive in the sample size $m$ as opposed to the dimensionality of the configuration space. Note that $\int d\vec{W} \rho_m(\vec{W}) = 1$, and that the partition function $Z_m$ provides the normalization factor. Note also that this distribution arises without invoking specific algorithms for using $E_L(\vec{W})$ to explore the configuration space $\{\vec{W}\}$.

The training data $\vec{\xi} = (\vec{x}, \vec{y})$ is drawn from a natural distribution

$$\tilde{P}(\vec{\xi}) = \tilde{P}(\vec{x}, \vec{y}) = \tilde{P}(\vec{y} \mid \vec{x}) \tilde{P}(\vec{x}). \tag{24}$$

Here, $\tilde{P}(\vec{x})$ describes the region of interest input space, and $\tilde{P}(\vec{y} \mid \vec{x})$ describes the functional dependence. The partition function

$$Z_m = \int d\vec{W} \rho_0(\vec{W}) \exp\left(-\beta \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{\xi}^\mu\right)\right) \tag{25}$$

depends on the specific set of data points $D = \{\vec{\xi}^\mu\}$ drawn from $\tilde{P}(\vec{\xi})$. The associated free energy

$$F = -(1/\beta) \langle\langle \ln Z_m \rangle\rangle_D \tag{26}$$

follows from averaging over all possible data sets of size $m$. The average learning error follows from the usual thermodynamic derivative:

$$E_L = -\frac{1}{m} \frac{\partial}{\partial \beta} \langle\langle \ln Z_m \rangle\rangle_D \tag{27}$$

The entropy follows from

$$F = m E_L - (1/\beta) S \tag{28}$$

For the learning process, this results in

$$S = -\int d\vec{W} \rho_m(\vec{W}) \ln\left[\frac{\rho_m(\vec{W})}{\rho_0(\vec{W})}\right] = -D_{KL}\left[\rho_m \mid \rho_0\right]. \tag{29}$$
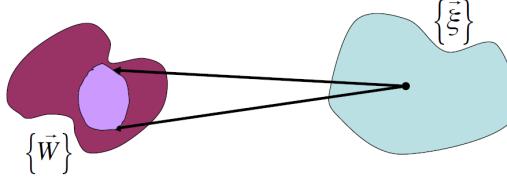
FIG. 8. Relation between the configuration space of the model and the sample space. $P(\vec{W}) = \rho_0(\vec{W})$ : prior distribution (left, dark purple). $P(\vec{W} \mid \vec{\xi})$ : distribution of hypothesis induced by example $\vec{\xi}$ (left, light purple).

The entropy of learning is minus the Kullback-Leibler distance between the posterior $\rho_m(\vec{W})$ and the prior $\rho_0(W)$; this distance measures the amount of information gained. The distance between posterior and prior increases monotonically with the size $m$ of the training set.

The entropy difference $\Delta H = H_{P(\vec{W})} - \left\langle\left\langle H_{P(\vec{W}|\vec{\xi})} \right\rangle\right\rangle_{P(\vec{\xi})}$ equals the mutual information between the $\{\vec{W}\}$ space (the model) and the $\{\vec{\xi}\}$ space (the world) (Figs. 8 and 9).

### D. Maximum likelihood learning

There are two approaches to learning: minimizing the error on the data,

$$E_L(\vec{W}) = \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{\xi}^{\mu}\right), \tag{30}$$

and maximizing the likelihood of the data,

$$\mathcal{L}(\vec{W}) = \prod_{\mu=1}^{m} P\left(\vec{\xi}^{\mu} \mid \vec{W}\right). \tag{31}$$



FIG. 9. Relation between the true sample distribution (right, light blue) and the distribution induced by the choice of hypothesis $\vec{W}$ (right, light purple).

The likelihood of the data given the network model can be expressed as

$$\mathcal{L}(\vec{W}) = P(D \mid \vec{W}) = P\left(\vec{\xi}^1, \vec{\xi}^2, \ldots, \vec{\xi}^m \mid \vec{W}\right) = \prod_{\mu=1}^{m} P\left(\vec{\xi}^{\mu} \mid \vec{W}\right) \tag{32}$$

But, what is the form of $P(\vec{\xi} \mid \vec{W})$? We require that these two approaches be coherent: that the minimization of the learning error implies the maximization of the data likelihood, and viceversa. This requirement leads to an explicit form for the likelihood:

$$P(\vec{\xi} \mid \vec{W}) = \frac{1}{z(\beta)} \exp(-\beta E(\vec{W} \mid \vec{\xi})), \tag{33}$$

where $z(\beta)$ is a normalization factor that guarantees that $\int P(\vec{\xi} \mid \vec{W}) \, d\vec{\xi} = 1$.

We can now compute the likelihood of the data:

$$P(D \mid \vec{W}) = \prod_{\mu=1}^{m} P\left(\vec{\xi}^{\mu} \mid \vec{W}\right) \tag{34}$$

$$= \frac{1}{z(\beta)^m} \exp\left(-\beta \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{\xi}^{\mu}\right)\right) \tag{35}$$

$$= \frac{1}{z(\beta)^m} \exp\left(-\beta m E_L(\vec{W})\right) . \tag{36}$$

Bayesian inference

$$P(\vec{W} \mid D) = \frac{P(D \mid \vec{W}) \, P(\vec{W})}{P(D)} \tag{37}$$

then leads to a the Gibbs distribution

$$\rho_m(\vec{W}) = \frac{1}{Z_m} \rho_0(\vec{W}) \exp\left(-\beta m E_L(\vec{W})\right) . \tag{38}$$

This is now consistent with the exponential decay of survival probability we chose before.

We can now check for the correspondence between the two formulations:

| | **Bayes** | $\leftrightarrow$ | **Gibbs** |
|---|---|---|---|
| Prior: | $P(\vec{W})$ | $\leftrightarrow$ | $\rho_0(\vec{W})$ |
| Posterior: | $P(\vec{W} \mid D)$ | $\leftrightarrow$ | $\rho_m(\vec{W})$ |
| Likelihood: | $P(D \mid \vec{W})$ | $\leftrightarrow$ | $\frac{1}{z(\beta)^m} \exp\left(-\beta m E_L(\vec{W})\right)$ |
| Evidence: | $P(D)$ | $\leftrightarrow$ | $\frac{1}{z(\beta)^m} Z_m$ |

where $P(D) = \int d\vec{W} P(D \mid \vec{W}) P(\vec{W})$.

### 1. Generalization ability

The normalization constant $z(\beta)$ plays a role in the evaluation of prediction errors (*i.e.,* has the network model acquired a good model of the world?).

Consider now a new point $\vec{\xi}$ not part of the training data $D = \left\{\vec{\xi}^1, \vec{\xi}^2, \ldots, \vec{\xi}^m\right\}$. What is the likelihood of this test point?

$$P(\vec{\xi} \mid D) = \int d\vec{W} P(\vec{\xi} \mid \vec{W}) \, P(\vec{W} \mid D) , \tag{39}$$

with

$$P(\vec{\xi} \mid \vec{W}) = \frac{1}{z(\beta)} \exp(-\beta E(\vec{W} \mid \vec{\xi})) \tag{40}$$

and

$$P(\vec{W} \mid D) = \rho_m(\vec{W}) = \frac{1}{Z_m} \rho_0(\vec{W}) \exp\left(-\beta \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{\xi}^{\mu}\right)\right) . \tag{41}$$

Thus,

$$P(\vec{\xi} \mid D) = \int d\vec{W} \; P(\vec{\xi} \mid \vec{W}) \, P(\vec{W} \mid D) \tag{42}$$

$$= \frac{1}{z(\beta) Z_m} \int d\vec{W} \rho_0(\vec{W}) \exp\left(-\beta \sum_{\mu=1}^{m+1} E\left(\vec{W} \mid \vec{\xi}^{\mu}\right)\right) , \tag{43}$$

where $\vec{\xi}^{m+1} = \vec{\xi}$ is the test point; it appears as if it had been added to the training set. Then

$$P(\vec{\xi} \mid D) = \frac{Z_{m+1}}{z(\beta)Z_m} \,. \tag{44}$$

The generalization error of the trained model is defined as the logarithm of the likelihood of an arbitrary test point $\vec{\xi}$ drawn from the same distribution as the training data:

$$P(\vec{\xi} \mid D) = \frac{Z_{m+1}}{z(\beta)Z_m} \implies E_G = -\frac{1}{\beta}\left[\ln\frac{Z_{m+1}}{Z_m} - \ln z(\beta)\right] \,. \tag{45}$$

For large $m$, the difference between $\ln Z_{m+1}$ and $\ln Z_m$ can be approximated by a derivative with respect to $m$. Then $\ln Z$ is averaged over all possible data sets of size $m$, to obtain:

$$E_G = -\frac{1}{\beta}\frac{\partial}{\partial m}\langle\langle \ln Z_m\rangle\rangle_D + \frac{1}{\beta}\ln z(\beta) \,. \tag{46}$$

This additional thermodynamic derivative is unusual within a canonical Gibbs ensemble. It arises because the role of the energy in the exponential factor is played by the learning error; this quantity is extensive in the number $m$ of examples in the training set, which differs from the dimensionality $D_{\vec{W}}$ of the configuration space. It is thus possible to take a derivative with respect to $m$ at constant $D_{\vec{W}}$.

Learning and generalization errors thus correspond to the two thermodynamic derivatives:

$$E_L = -\frac{1}{m}\frac{\partial}{\partial\beta}\langle\langle \ln Z_m\rangle\rangle_D \tag{47}$$

$$E_G = -\frac{1}{\beta}\frac{\partial}{\partial m}\langle\langle \ln Z_m\rangle\rangle_D + \frac{1}{\beta}\ln z(\beta) \,. \tag{48}$$

What remains to be determined is how to select the tolerance to errors, the inverse temperature $\beta$. In the following, we will show how this arises in a simple learning scenario.

### E.   A simple example: The linear map
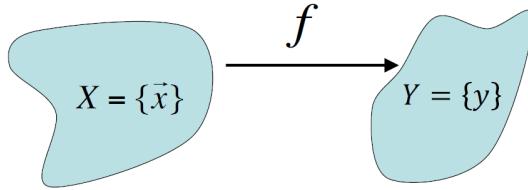


FIG. 10. Linear mapping $f : X \to Y$.

We consider the simplest possible network: a linear map from an $N$-dimensional input $\vec{x}$ to an scalar $y$. The parameters $\{\vec{W}\}$ are drawn from $\rho_0(\vec{W}) = \mathcal{N}(0, C_w)$ with $C_w = \sigma_w^2 I_N$ and $\sigma_w \gg 1$. The map is described by

$$\vec{x} = \{x_1, x_2, \ldots, x_N\} \to y = f_W(\vec{x}) = \sum_{i=1}^{N} W_i x_i = \vec{W}^T\vec{x} \,. \tag{49}$$

Let's assume that the inputs are also drawn from a Gaussian distribution $\tilde{P}(\vec{x}) = \mathcal{N}(0, C_x)$ with $C_x = \sigma_x^2 I_N$, and that the target output for input $\vec{x}^\mu$ is

$$y^\mu = \vec{W}_0^T\vec{x}^\mu + \eta^\mu \tag{50}$$

such that

$$\tilde{P}(y \mid \vec{x}) = \mathcal{N}\left(\vec{W}_0^T\vec{x}^\mu, \sigma_\eta^2\right) \,. \tag{51}$$

We then train to minimize

$$E_L(\vec{W}) = \frac{1}{2}\sum_{\mu=1}^{m}\left(y^\mu - \vec{W}^T\vec{x}^\mu\right)^2 = \frac{1}{2}\sum_{\mu=1}^{m}\left(\left(\vec{W} - \vec{W}_0\right)^T\vec{x}^\mu - \eta^\mu\right)^2 \tag{52}$$

For $\sigma_w \gg 1$ and in the large $m$ limit the free energy can be computed analytically [4]:

$$\langle\langle\ln Z_m\rangle\rangle = -N\ln\sigma_w - \frac{\mathbf{W}_0^T\cdot\mathbf{W}_0}{2\sigma_w^2} - \frac{N}{2}\ln\left(2\beta\sigma_x^2 m\right) \tag{53}$$

$$+ (N-m)\beta\sigma_\eta^2 + O(1/m)\,. \tag{54}$$

The thermodynamic derivatives are:

$$E_L = \frac{N}{2m\beta} + \left(1 - \frac{N}{m}\right)\sigma_\eta^2 + O\left(1/m^2\right) \tag{55}$$

$$E_G = \frac{N}{2m} + \beta\sigma_\eta^2 + \ln\left[\frac{\pi}{\beta}\right]^{1/2} + O\left(1/m^2\right)\,. \tag{56}$$

We can now define the effective temperature $\beta_0$ associated to the noise in the examples as

$$\beta_0 = \frac{1}{2\sigma_\eta^2}\,. \tag{57}$$

The thermodynamic derivatives then simplify as

$$E_L = \frac{N}{2m}\left[\frac{1}{\beta} - \frac{1}{\beta_0}\right] + \frac{1}{2\beta_0} + O\left(1/m^2\right)$$

$$E_G = \frac{N}{2m} + \frac{\beta}{2\beta_0} + \ln\left[\frac{\pi}{\beta}\right]^{1/2} + O\left(1/m^2\right)\,. \tag{58}$$
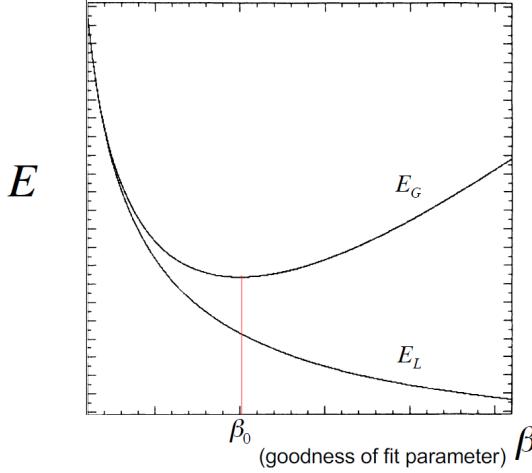


FIG. 11. Learning error $E_L$ and generalization error $E_G$ as function of the inverse temperature $\beta$.

For small values of $\beta$, in the high temperature regime, tolerance to learning error is large. The data is being underfit, and both $E_L$ and $E_G$ are large. As $\beta$ increases, tolerance to learning error decreases and $E_L$ decreases monotonically. While $E_G$ exceeds $E_L$ for all values of $\beta$, as expected, a different regime in $E_G$ arises for $\beta > \beta_0$. In this large $\beta$, low temperature regime, the tolerance to learning error continues to decrease but $E_G$ increases, indicating that the data is overfit. This nonmonotonic behavior in $E_G$ is observed in a variety of overfitting scenarios. Here is arises when the effective temperature of the learning algorithm falls below the value that characterizes the level of noise in the data.

### F.  Appendix: Derivation of the exponential form for the prediction error

Let's require that the minimization of the learning error

$$E_L(\vec{W}) = \sum_{\mu=1}^{m} E\left(\vec{W} \mid \vec{\xi}^{\mu}\right) \tag{59}$$

guarantees the maximization of the likelihood

$$\mathcal{L}(\vec{W}) = \prod_{\mu=1}^{m} P\left(\vec{\xi}^{\mu} \mid \vec{W}\right) . \tag{60}$$

Given a training set $\left(\vec{\xi}^1, \vec{\xi}^2, \ldots, \vec{\xi}^m\right)$, these two functions need to be related:

$$\mathcal{L}(\vec{W}) = \Phi\left(E_L(\vec{W})\right) . \tag{61}$$

Take a derivative on both sides with respect to one of the points in the training set, $\vec{\xi}^j$:

$$\frac{\partial \mathcal{L}(D \mid \vec{W})}{\partial \vec{\xi}^j} = \mathcal{L}\left(D \mid \vec{W}\right) \frac{1}{P\left(\vec{\xi}^j \mid \vec{W}\right)} \frac{\partial P\left(\vec{\xi}^j \mid \vec{W}\right)}{\partial \vec{\xi}^j} \tag{62}$$

$$= \Phi' \frac{\partial E\left(\vec{W} \mid \vec{\xi}^j\right)}{\partial \vec{\xi}^j} . \tag{63}$$

This leads to

$$\frac{\Phi'}{\Phi} = \frac{\frac{1}{P\left(\vec{\xi}^j \mid \vec{W}\right)} \frac{\partial P\left(\vec{\xi}^j \mid \vec{W}\right)}{\partial \vec{\xi}^j}}{\frac{\partial E\left(\vec{W} \mid \vec{\xi}^j\right)}{\partial \vec{\xi}^j}} . \tag{64}$$

While the left-hand side of the equation depends on the full training set $\left(\vec{\xi}^1, \vec{\xi}^2, \ldots, \vec{\xi}^m\right)$, the right-hand side depends only on $\vec{\xi}^j$. The only way for this equality to hold for all values of $\left(\vec{\xi}^1, \vec{\xi}^2, \ldots, \vec{\xi}^m\right)$ is for both sides to be actually independent of the data, and thus equal to a constant:

$$\frac{\frac{1}{P\left(\vec{\xi}^j \mid \vec{W}\right)} \frac{\partial P\left(\vec{\xi}^j \mid \vec{W}\right)}{\partial \vec{\xi}^j}}{\frac{\partial E\left(\vec{W} \mid \vec{\xi}^j\right)}{\partial \vec{\xi}^j}} = -\beta . \tag{65}$$

The equation

$$\frac{1}{P\left(\vec{\xi}^j \mid \vec{W}\right)} \frac{\partial P\left(\vec{\xi}^j \mid \vec{W}\right)}{\partial \vec{\xi}^j} = -\beta \frac{\partial E\left(\vec{W} \mid \vec{\xi}^j\right)}{\partial \vec{\xi}^j} \tag{66}$$

can be integrated to obtain

$$P\left(\vec{\xi}^j \mid \vec{W}\right) \propto \exp\left(-\beta E\left(\vec{W} \mid \vec{\xi}^j\right)\right) . \tag{67}$$

The normalized probability distribution is

$$P(\vec{\xi} \mid \vec{W}) = \frac{1}{z(\beta)} \exp(-\beta E(\vec{W} \mid \vec{\xi})) , \tag{68}$$

with $z(\beta) = \int d\vec{\xi} \exp(-\beta E(\vec{W} \mid \vec{\xi}))$. Since the equation that determines $P(\vec{\xi} \mid \vec{W})$ is first order, there is only one constant of integration: $\beta$. For $\beta > 0$, minima of $E$ correspond to maxima of $P$, as required.

## G.  Summary

- A Gibbs probability density function describes the ensemble of trained networks; it corresponds to the Bayesian posterior.

- The normalization of this probability density corresponds to the Gibbs partition function, from which a free energy can be obtained.

- The thermodynamic formulation provides the learning error and a relative entropy that quantifies the exchange of information between the model and the world as a provider of examples to be learned.

- A correspondence between error minimization and likelihood maximization provides an equation for the generalization error that leads to a novel thermodynamic derivative.

## II. LECTURE 2: GENERALIZED LINEAR MODELS *VS.* BACK PROPAGATION THROUGH TIME

We will now focus on neural networks in the biological context, and on how to analyze neural population dynamics to find lower dimensional representations. We start with the simultaneous recording of the activity of a population of neurons in the brain of a nonhuman primate. High-density multi-electrode arrays (MEAs) implanted into the cortex allow for the simultaneous recording of the order of a hundred neurons (Fig. 12). Given this data, we aim to characterize the dynamics of the underlying neural circuits.
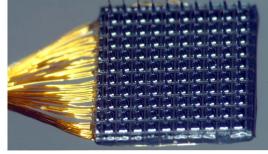


FIG. 12. Utah array for simutaneously recording the activity of multiple neurons.

### A. Neural population activity

Population recordings provide the temporal pattern of neuron spikes, the spike trains (Fig. 13). To visualize these spike trains we use one row per neuron and display a dot whenever that neuron spikes. The vertical axis labels the recorded neurons, the horizontal axis is time. Consider a population of $N$ neurons whose spiking activity is observed during a time interval $(0, T]$. The interval is divided into $K$ bins of size $\Delta = T/K$, labeled by an index $1 \leq k \leq K$. In each interval $k$ we observe the number of spikes $y_i(k)$ emitted by neuron $i$, for all $1 \leq i \leq N$. The number and distribution of spikes varies over neurons, over time, and depends on the settings of the experiment.
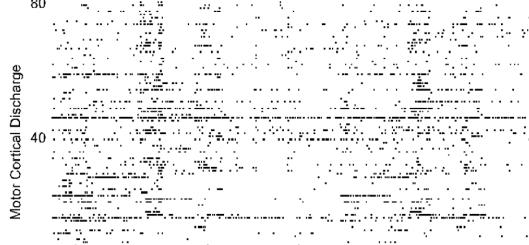


FIG. 13. Neural spiking activity for a population of neurons in primary motor cortex.
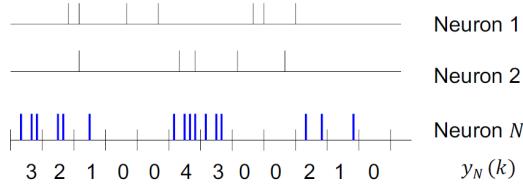


FIG. 14. Binned spike train.

If the subject is performing a task during the interval $(0, T]$, we search for task-specific patterns in the data. Consider a simple motor task: *center-out reaches*. The subject is instructed to reach a visually displayed target from a center towards the outside, with an instructed delay after the target is shown (Fig. 15). During the execution of this task, some neurons in the primary motor cortex M1 exhibit a firing rate that is modulated by the direction of the reach.
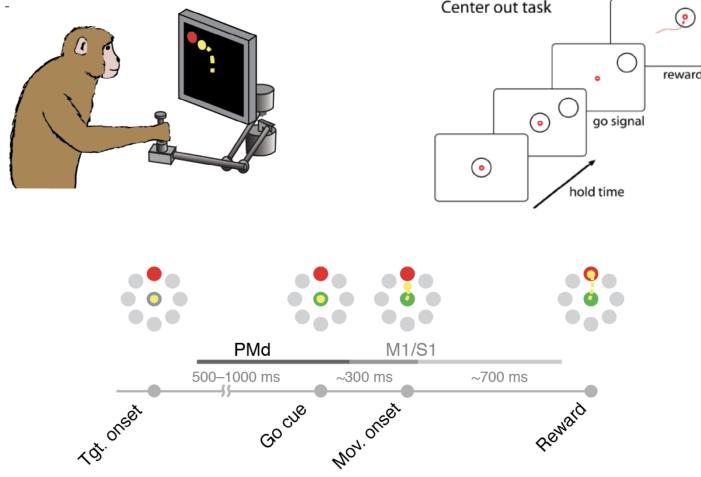
FIG. 15. Instructed delay center-out reaching task

### 1. Analysis of neural population activity

In our analysis we will assume that the distribution of the spike counts $y$ within a bin can be described by a Poisson distribution:

$$\rho(y \mid \lambda) = \frac{\lambda^y e^{-\lambda}}{y!} \tag{69}$$

It can be easily checked that the distribution is properly normalized

$$\sum_{y=0}^{\infty} \frac{\lambda^y e^{-\lambda}}{y!} = e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} = e^{-\lambda} e^{+\lambda} = 1 \tag{70}$$

and has moments

$$\mathrm{E}(y) = \langle y \rangle = \lambda \tag{71}$$

$$\mathrm{Var}(y) = \left\langle (y - \langle y \rangle)^2 \right\rangle = \lambda \tag{72}$$

and Fano factor

$$\frac{\mathrm{Var}(y)}{\mathrm{E}(y)} = 1 \ . \tag{73}$$

Recordings from individual M1 neurons during the execution of a center-out task reveal both trial-to-trial variability and specific firing patterns that are direction dependent (Fig. 16). How can we model both variability and specificity? For Poisson statistics, the parameter $\lambda_i(k)$ is the mean or expectation value of the random variable $y_i(k)$. The trial-to-trial fluctuations of $y_i(k)$ about its mean $\lambda_i(k)$ describe the variability of neural activity within the $k$th time bin. The time-dependent parameter $\lambda_i(k)$ provides a tool for specificity: we will model $\lambda_i(k)$ through its relation to sensory stimuli, motor output, and the spiking activity of other neurons.

### 2. Neural activity as a Poisson process

Let's go back to the population data $\{y_i(k)\}$, for $1 \leq k \leq K$ and $1 \leq i \leq N$. The spiking activity of neuron $i$ at time interval $k$ is modeled as a Poisson process with mean $\lambda_i(k)$. The probability of observing precisely $y_i(k)$ spikes emitted by neuron $i$ at time $k$ is given by:

$$P(y_i(k) \mid \lambda_i(k)) = \frac{(\lambda_i(k))^{y_i(k)} e^{-\lambda_i(k)}}{y_i(k)!} \ . \tag{74}$$
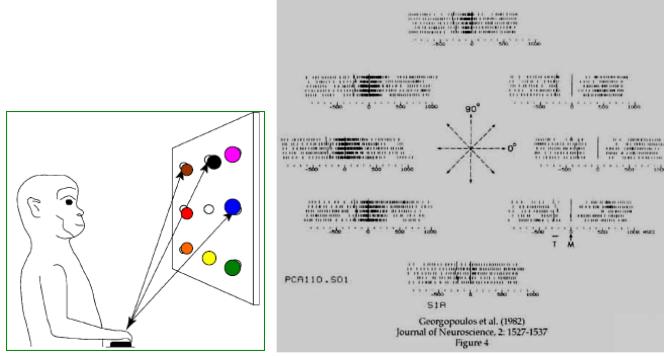
FIG. 16. Variability of M1 neural responses during a center-out reaching task.

What is then the probability of the observed data $\{y_i(k)\}$ given the parameters $\{\lambda_i(k)\}$? Within the Poisson assumption, the conditional probability is given by

$$P_D\left(\{y_i(k)\} \mid \{\lambda_i(k)\}\right) = \prod_{i=1}^{N}\prod_{k=1}^{K} \frac{(\lambda_i(k))^{y_i(k)} e^{-\lambda_i(k)}}{y_i(k)!} \ . \tag{75}$$

The log-likelihood, the logarithm of the probability, is then given by

$$L_D\left(\{y_i(k)\} \mid \{\lambda_i(k)\}\right) = \ln\left\{\prod_{i=1}^{N}\prod_{k=1}^{K} \frac{(\lambda_i(k))^{y_i(k)} e^{-\lambda_i(k)}}{y_i(k)!}\right\} \tag{76}$$

$$= \left\{\sum_{i=1}^{N}\sum_{k=1}^{K} [y_i(k)\ln\lambda_i(k) - \lambda_i(k) - \ln(y_i(k)!)]\right\} \ . \tag{77}$$

Here, $\{y_i(k)\}$ is the data while the neuron specific and time specific firing rates $\{\lambda_i(k)\}$ are the parameters of the model. If we were to estimate the parameters $\{\lambda_i(k)\}$ that maximize the likelihood of the data $\{y_i(k)\}$, what would we obtain? The answer is $\lambda_i(k) = <y_i(k)>$. This answer, based on the empirical mean of the data, is not satisfactory if our goal is to model the $\{\lambda_i(k)\}$.

### B. A model for the time dependent and neurons specific parameters $\{\lambda_i(k)\}$

Our goal is to model the parameters $\{\lambda_i(k)\}$ in order to relate their values to sensory stimuli, motor outputs, and the activity of other neurons in the network. We will approach this modeling task with generalized linear models (GLMs). But first we will make a detour into statistics: the exponential family of probability distributions.

#### 1. The exponential family of probability distributions

The exponential family encompasses probability distributions of the form:

$$\rho(y \mid \delta, \varphi) = \exp\left\{\frac{y\delta - b(\delta)}{a(\varphi)} + c(y, \varphi)\right\} \ . \tag{78}$$

Here, $y$ is the random variable whose probability density function is given by $\rho$. The distribution is parametrized by the canonical parameter $\delta$ and the dispersion parameter $\varphi$. The functions $a$, $b$, and $c$ need to be specified; they define the various distributions within the family.

The term $c(y, \varphi)$ plays an important role: it provides a normalization function that guarantees $\int dy\, \rho(y \mid \delta, \varphi) = 1$ for all $\delta, \varphi$. Since $\int dy\, \rho_y(y \mid \delta, \varphi) = 1$ for all $\delta, \varphi$ then:

$$\begin{aligned}
\frac{\partial}{\partial\delta}\int dy\rho_y(y \mid \delta, \varphi) = 0 &\Rightarrow \mathrm{E}(y) = b'(\delta) \\
\frac{\partial^2}{\partial\delta^2}\int dy\rho_y(y \mid \delta, \varphi) = 0 &\Rightarrow \mathrm{Var}(y) = a(\varphi)b''(\delta) \ .
\end{aligned} \tag{79}$$

Note that the canonical parameter $\delta$ fully determines the mean $\mathrm{E}(y)$ through $b(\delta)$, while the variance $\mathrm{Var}(y)$ requires additional information provided by the dispersion parameter through $a(\varphi)$ [5].

Consider the family of canonical exponential distributions with canonical parameter $\delta$ and dispersion parameter $\varphi$. Why is this family important? Because the normal, Bernoulli, binomial, multinomial, Poisson, gamma, geometric, chi-square, beta, and a few other distributions are all members of this exponential family.

### 2.  *The Poisson distribution as a member of the exponential family*

Let's have a closer look at one specific example, the Poisson distribution:

$$\rho(y \mid \lambda) = \frac{\lambda^y e^{-\lambda}}{y!} = \exp\{y \ln \lambda - \lambda - \ln(y!)\} . \tag{80}$$

This takes the form of the exponential family (Eq. 78) for $a(\varphi) = 1, \delta = \ln \lambda, b(\delta) = \lambda$, and $c(\varphi, y) = -\ln(y!)$.

The relations $\mathrm{E}(y) = b'(\delta)$ and $\mathrm{Var}(y) = a(\varphi)b''(\delta)$ hold for any probability density function within the exponential family. When applied to the Poisson case, they imply:

$$\begin{aligned} \mathrm{E}(y) &= b'(\delta) = \lambda \\ \mathrm{Var}(y) &= a(\varphi)b''(\delta) = \lambda . \end{aligned} \tag{81}$$

For Poisson statistics, $\mathrm{E}(y) = \mathrm{Var}(y) = \lambda$ implies:

$$b'(\delta) = a(\varphi)b''(\delta) \implies \begin{cases} a(\varphi) = 1 \\ b(\delta) = e^\delta = \lambda \end{cases} . \tag{82}$$

In a generalized linear model (GLM) for a probability distribution that is a member of the exponential family [5], the expectation value $\mathrm{E}(y)$ is related to the canonical parameter $\delta$ via a nonlinear link function $g$:

$$g(\mathrm{E}(y)) = \delta \qquad \mathrm{E}(y) = g^{-1}(\delta) . \tag{83}$$

This is the only nonlinearity in the model, as the canonical parameter $\delta$ is constructed as a linear combination of all observed variables that can *explain* the random variable $y$. In the Poisson case, $\delta = \ln \lambda = \ln(\mathrm{E}(y))$, and the nonlinear link function $g$ is the logarithm:

$$\begin{aligned} \lambda &= \mathrm{E}(y) = g^{-1}(\delta) = \exp(\delta) \\ \delta &= g(\lambda) = \ln(\lambda) \end{aligned} \tag{84}$$

### C.  Generalized linear models for spiking neurons

The parameter $\lambda_i(k)$ is the time-dependent and neuron-specific mean of a Poisson process. In a GLM for a Poisson distribution, it is the logarithm (link function) of the mean that is expressed as a linear combination of all observed variables that can be used to explain the observed firing rates. Here, we have

- Internal covariates: preceding neural activity (hidden neurons)

- External covariates: sensory stimulus (input), direction of motion (output)

We know the spiking history of the ensemble of $N$ neurons up to the current time $t$. We denote this as the spiking history of the ensemble:

$$H(t) = \left\{ \{y_i(t')\}_{i=1}^N, t' \le t \right\} \tag{85}$$

Given this information, what is our expectation of the number of spikes that neuron $i$ will fire in the interval $(t, t + \Delta)$? This is the *conditional intensity* $\lambda_i(t \mid H(t))$, a strictly positive function that provides a history-dependent generalization of the time dependent rate of an inhomogeneous Poisson process [6].
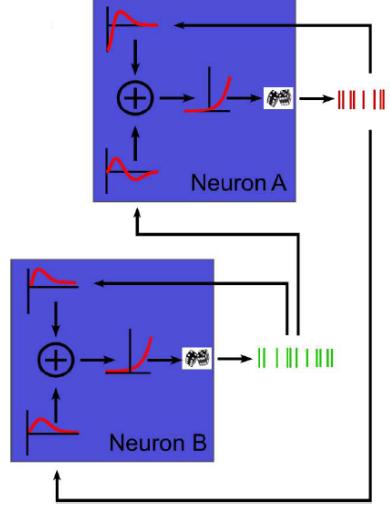
FIG. 17. History dependent GLM model for the time-dependent neuron-specific mean firing rate.

### 1. A GLM based on internal covariates

The history-dependent model for $\lambda_i(t \mid H(t))$ in the generalized linear model (GLM) approach is given by:

$$\delta_i(t \mid H(t), \{\alpha\}) = \ln \lambda_i(t \mid H(t), \{\alpha\}) = \alpha_{i0} + \sum_{j=1}^{N} \sum_{m=1}^{\tau_N} \alpha_{ij}(m) y_j(t-m) , \qquad (86)$$

which leads to a linear-nonlinear model for $\lambda_i(t)$:

$$\lambda_i(t \mid H(t), \{\alpha\}) = \exp\left\{ \alpha_{i0} + \sum_{j=1}^{N} \sum_{m=1}^{\tau_N} \alpha_{ij}(m) y_j(t-m) \right\} . \qquad (87)$$

Once the model model parameters $\{\alpha_{ij}(m)\}$ have been specified, this equation for $\lambda_i(t)$ provides a generative model (Fig. 17). The kernel parameter $\alpha_{ij}(m)$ quantifies the effect that the spiking activity of neuron $j$ at time bin $(t-m)$ has on the spiking activity of neuron $i$ at time bin $t$. For a given $(i, j)$ pair, the time dependent $\alpha_{ij}(m)$ defines an effective connectivity (Fig. 18).

How are the model parameters $\{\alpha_{ij}(m)\}$ determined? Given the data $\{y_i(k)\}$, we search for parameters that maximize the likelihood of the data:

$$L_D\left(\{y_i(k)\}\right) \propto \left\{ \sum_{i=1}^{N} \sum_{k=1}^{K} (y_i(k) \ln \lambda_i(k) - \lambda_i(k)) \right\} . \qquad (88)$$

A term that does not depend on $\{\lambda_i(k)\}$ and thus does not depend on the parameters has been dropped. Explicitly, the likelihood function to be maximized is:

$$L_D\left(\{y_i(k)\} \mid \{\alpha\}\right) \propto \left\{ \sum_{i=1}^{N} \sum_{k=1}^{K} \left( y_i(k) \left[ \alpha_{i0} + \sum_{j=1}^{N} \sum_{m=1}^{\tau_v} \alpha_{ij}(m) y_j(k-m) \right] \right. \right.$$

$$\left. \left. - \exp\left[ \alpha_{i0} + \sum_{j=1}^{N} \sum_{m=1}^{\tau_N} \alpha_{ij}(m) y_j(k-m) \right] \right) \right\} . \qquad (89)$$

### 2. Learning a history dependent GLM by gradient ascent

The maximization of the likelihood proceeds via gradient ascent with an adaptive step size. To implement this algorithm, we need to compute the first and second derivatives of the likelihood with respect to the parameters $\{\alpha_{ij}(m)\}$. The gradient
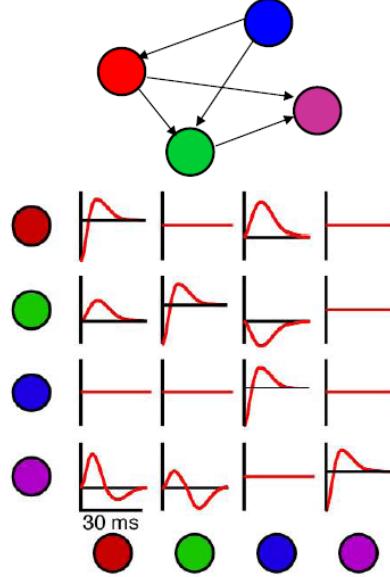
FIG. 18. Effective connectivity.

that drives the uphill search is given by:

$$\left.\frac{\partial L_D\left(\{y_i(k)\} \mid \{\alpha\}\right)}{\partial \alpha_{ij}(m)}\right|_{(\mu)} = \sum_{k=1}^{K}\left\{\left[y_i(k) - \langle y_i(k)\rangle^{(\mu)}\right] y_j(k-m)\right\} \tag{90}$$

The update of the parameter $\alpha_{ij}(m)$ is given by the product of the activity $y_j(k-m)$ of the *presynaptic* neuron $j$ at time lag $m$ and the difference between the actual activity $y_i(k)$ of the *postsynaptic* neuron and our current estimate of it at iteration $(\mu)$. The rule *presynaptic activity times postsynaptic error* is a famous learning rule, called the Delta Rule [1]. We italicize *presynaptic* and *postsynaptic* to avoid the implication that the parameter $\alpha_{ij}(m)$ is an actual synaptic strength.

The components of the Hessian matrix of second derivatives that controls the size of the uphill steps are given by:

$$\left.\frac{\partial^2 L_D\left(\{y_i(k)\} \mid \{\alpha\}\right)}{\partial \alpha_{ij}(m)\,\partial \alpha_{ij'}(m')}\right|_{(\omega)} = -\sum_{k=1}^{K}\langle y_i(k)\rangle^{(\mu)} y_j(k-m) y_{j'}(k-m'). \tag{91}$$

Now there are two *presynaptic* neurons: neuron $j$ at time lag $m$ and neuron $j'$ at time lag $m'$. Their activities are multiplied, and this product is weighted by our current estimate of the activity of the *postsynaptic* neuron $i$. Note the overall minus sign! The variables $\{y\}$ represent number of spikes emitted during a bin of size $\Delta$. These variables and their averages are always non-negative. Every component of the Hessian matrix is negative - the surface is everywhere convex!

The gradient ascent algorithm can be written as follows:

$$\vec{\alpha}^{(\mu+1)} = \vec{\alpha}^{(\mu)} + \mathrm{E}^{(\mu)}\vec{\nabla}L^{(\mu)} \tag{92}$$

Here, $\vec{\alpha}$ is a listing of all the parameters needed to specify the model; $\vec{\nabla}L$ is the gradient of the likelihood function $L$, obtained by taking the derivative of $L$ with respect to each parameter in $\vec{\alpha}$; and E is the matrix of step sizes, obtained by inverting the Hessian matrix of second derivatives of the likelihood function. If the model requires $p$ parameters, then both $\vec{\alpha}$ and $\vec{\nabla}L$ are $p$-dimensional vectors, and E is a $p \times p$ matrix.

### D. An example: a GLM for the center-out task

We show the analysis of data acquired from two humans with tetraplegia while they participated in a clinical trial. The data consists of M1 recordings (Fig. 19) obtained while the participants performed the center-out task; the (x,y) motion of the cursor followed from decoding M1 neural activity [7].
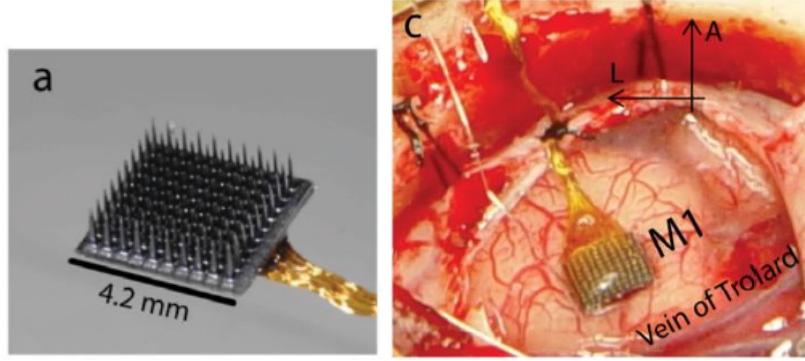
FIG. 19. A multi-electrode array (left) that allows for the simultaneous recording of about a hundred neurons is shown when inserted in primary motor cortex M1 (right).

### 1. A GLM based on spiking history for the center-out task

To formulate the GLM, we ask what is the probability that neuron $i$ spikes at bin $k$, conditioned on the spiking history $H_k$ of the neural population during the preceding 100 ms:

$$\ln \lambda_i (k \mid H_k) = \alpha_{i0} + \sum_{m=1}^{\tau_N} \alpha_{ii}(m) y_i(k-m) + \sum_{j=1, j \neq i}^{N} \sum_{m=1}^{\tau_N} \alpha_{ij}(m) y_j(k-m) . \tag{93}$$

Here, $\Delta = 1$ ms and $\tau_N = 100$. Data is used to fit the conditional intensity $\lambda_i(k)$ and obtain the background level $\alpha_{i0}$ of spiking activity for neuron $i$, the kernel $\alpha_{ii}(m)$ related to its intrinsic history effects, and the kernels $\alpha_{ij}(m)$ related to history effects due to the other neurons $j \neq i$. Once the model is fitted, the estimated probability of a spike at any time bin can be computed.

The fitting of the model parameters involves the following steps:

(1) Given the data $\{y_i(m)\}$ and the current value $\{\alpha^{(\mu)}\}$ of the parameters, construct:

$$\langle y_i(k) \rangle^{(\mu)} = \exp \left[ \alpha_{i0}^{(\mu)} + \sum_{j'=1}^{N} \sum_{m'=1}^{\tau_N} \alpha_{ij'}^{(\mu)}(m') y_{j'}(k-m') \right] \tag{94}$$

Once the estimates $\langle y_i(k) \rangle^{(\mu)}$ have been computed, the current values of the parameters are no longer needed.

(2) Build the components of the gradient vector:

$$\left. \frac{\partial L_D (\{y_i(k)\} \mid \{\alpha\})}{\partial \alpha_{ij}(m)} \right|_{(\mu)} = \sum_{k=1}^{K} \left[ y_i(k) - \langle y_i(k) \rangle^{(\mu)} \right] y_j(k-m) \tag{95}$$

(3) Build the components of the Hessian matrix:

$$\left. \frac{\partial^2 L_D (\{y_i(k)\} \mid \{\alpha\})}{\partial \alpha_{ij}(m) \partial \alpha_{ij'}(m')} \right|_{(\mu)} = -\sum_{k=1}^{K} \langle y_i(k) \rangle^{(\mu)} y_j(k-m) y_{j'}(k-m') \tag{96}$$
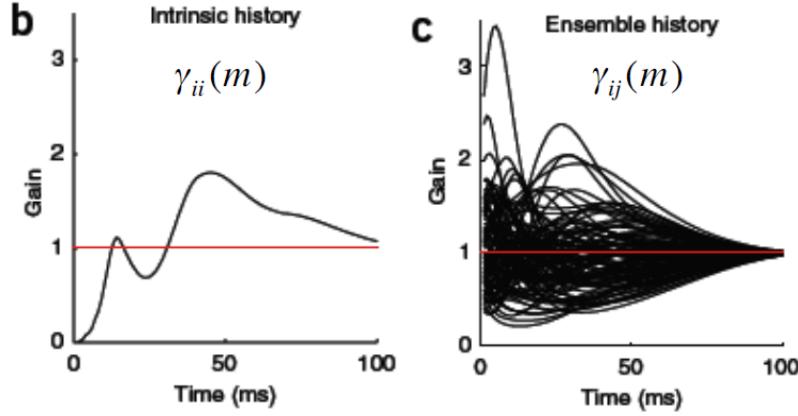
(4) Invert the Hessian matrix of second derivatives to obtain the matrix Epsilon of step sizes:

$$\mathrm{E} = -H^{-1} \tag{97}$$

(5) Multiply the matrix E and the gradient $\vec{\nabla} L$ to obtain the update:

$$\vec{\alpha}^{(\mu+1)} = \vec{\alpha}^{(\mu)} + \mathrm{E}^{(\mu)} \vec{\nabla} L^{(\mu)} \tag{98}$$

Examples of fitted temporal filters are shown for $i = 34$ using $\gamma_{ij}(m) = \exp \left( \alpha_{ij}(m) \right)$ (Fig. 20).

FIG. 20. $i = 34$          $i = 34, j \neq i$

*2. A GLM based on direction of motion for the center out task*

We now wish to use information about the direction of motion in order to predict the mean $\lambda_i(t)$ of the Poisson process. Information about a planar reach of extent $r$ in a direction characterized by $\theta$ can be extracted from the firing activity of orientation selective M1 neurons [8]:

$$f_i(r, \theta) = b_i + r a_i (1/2)[1 + \cos(\theta - \theta_i)] , \tag{99}$$

where $b_i$ is the background activity, $a_i$ is the amplitude of activity modulation, and $\theta_i$ is the preferred direction of neuron $i$ (Fig. 21).
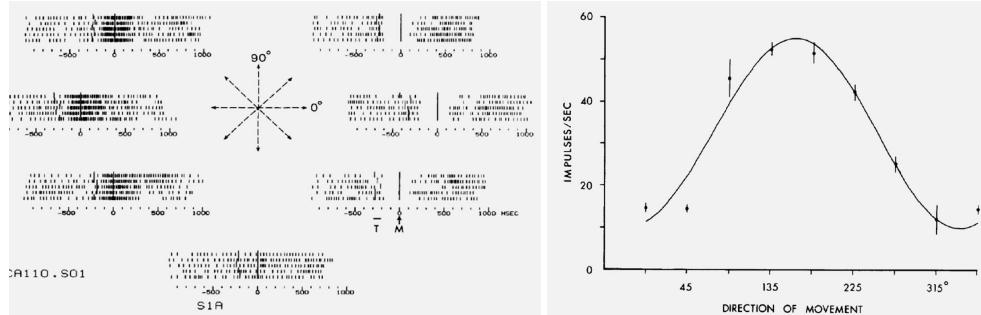


FIG. 21. The tuning curve is characterized by three model parameters. For experimental details, see [8].

This observation would allow to construct a hierarchy of three increasingly detailed GLM models:
(1) Non-interacting cosine-tuned neurons:

$$\ln \lambda_i(t) = \alpha_{i0} + (1/2)\alpha_{iR} r(t + \tau_R)[1 + \cos(\theta(t + \tau_R) - \theta_i)] \tag{100}$$

(2) History-dependent non-interacting cosine-tuned neurons:

$$\ln \lambda_i(t) = \alpha_{i0} + (1/2)\alpha_{iR} r(t + \tau_R)[1 + \cos(\theta(t + \tau_R) - \theta_i)] + \sum_{m=1}^{\tau_N} \alpha_{ii}(m) y_i(t - m) \tag{101}$$

(3) History-dependent interacting cosine-tuned neurons:

$$\ln \lambda_i(t) = \alpha_{i0} + (1/2)\alpha_{iR} r(t + \tau_R)[1 + \cos(\theta(t + \tau_R) - \theta_i)] + \sum_{j=1}^{N} \sum_{m=1}^{\tau_N} \alpha_{ij}(m) y_j(t - m) \tag{102}$$
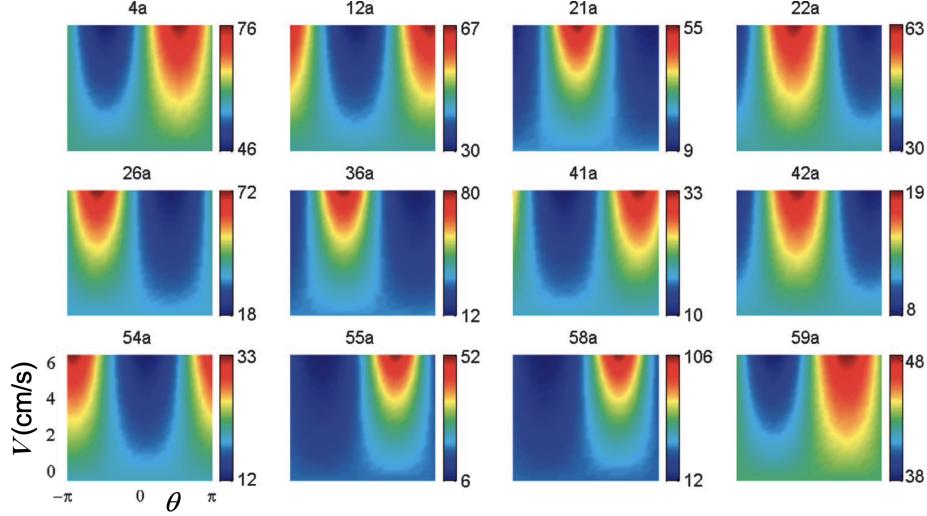
FIG. 22. Velocity tuning functions for 12 different neurons. The values of the mean number of spikes $\lambda$ are color coded.

### E. GLM models of increasing complexity

As an illustration of this approach, we review the analysis of MEA (MultiElectrode Array) recordings of neural activity in the arm area of primary motor cortex (M1) of awake and behaving monkeys involved in the execution of a two-dimensional tracking task: the pursuit and capture of a smoothly and randomly moving visual target [6]. The target was tracked by moving a two-link low friction manipulandum that constrained hand movementx to the horizontal plane. The $(x, y)$ hand position was digitized and resampled at 1 KHz; low-pass filtered finite differences of position were used to obtain the two components of velocity.

#### 1. Two GLM models for the tracking task

Two models were used to analyze the data:

(1) Velocity model [9]

(2) Velocity model plus intrinsic spiking history [6]

The velocity model uses $v_x(t + \tau_R)$ and $v_y(t + \tau_R)$ as explanatory variables for $\lambda_i(t)$:

$$
\begin{aligned}
&\lambda_i\left(t \mid v\left(t+\tau_R\right), \theta\left(t+\tau_R\right),\left\{\alpha_{i0}, \alpha_{iX}, \alpha_{iY}\right\}\right) \\
&= \exp\left\{\alpha_{i0} + \alpha_{iX} v_X\left(t+\tau_R\right) + \alpha_{iY} v_Y\left(t+\tau_R\right)\right\} \\
&= \exp\left\{\alpha_{i0} + v\left(t+\tau_R\right)\left[\alpha_{iX} \cos\left(\theta\left(t+\tau_R\right)\right) + \alpha_{iY} \sin\left(\theta\left(t+\tau_R\right)\right)\right]\right\}
\end{aligned}
\tag{103}
$$

Note that the model does not include a sum over time lags; it uses single time shift $\tau_R = 150$ ms. There are only three parameters to be determined through a maximum likelihood fit to the spiking data of the $i$th neuron: $\{\alpha_{i0}, \alpha_{iX}, \alpha_{iY}\}$. Once the values for these parameters have been specified, the conditional intensity $\lambda_i(t)$ can be plotted as a function of the subsequent velocity in polar coordinates: $v(t + \tau_R), \theta(t + \tau_R)$ (Fig. 22).

The next level of model complexity corresponds to including the intrinsic spiking history as an internal covariate:

$$
\begin{aligned}
&\lambda_i\left(t \mid H_i(t), v\left(t+\tau_R\right), \theta\left(t+\tau_R\right),\left\{\alpha_{i0},\left\{\alpha_{ii}(m)\right\}, \alpha_{iX}, \alpha_{iY}\right\}\right) \\
&= \exp\left\{\alpha_{i0} + v\left(t+\tau_R\right)\left[\alpha_{iX} \cos\left(\theta\left(t+\tau_R\right)\right) + \alpha_{iY} \sin\left(\theta\left(t+\tau_R\right)\right)\right] + \sum_{m=1}^{\tau_N} \alpha_{ii}(m) y_i(t-m)\right\} .
\end{aligned}
\tag{104}
$$

In addition to the three parameters $\{\alpha_{i0}, \alpha_{ix}, \alpha_{iY}\}$, this model for $\lambda_i(t)$ includes the parameters $\{\alpha_{ii}(m)\}$ that characterize the intrinsic history filter.
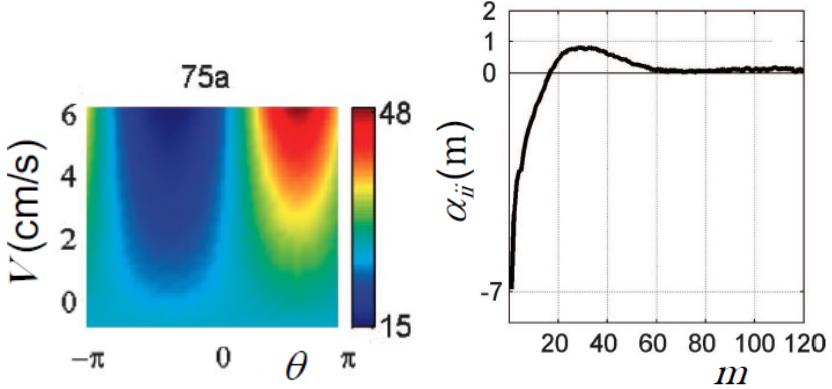
FIG. 23. Velocity tuning function (left) and intrinsic history filter (right) for neuron 75a.

This model was implemented by using a bin size $\Delta = 1$ ms. At this time resolution, the number of spikes $y_i(t-m)$ in any bin can only be 0 or 1. The maximum temporal length of the intrinsic history filter was set to $\tau_N = 120$. Results for neuron $i =$ 75a (Fig. 23) show distinctive velocity tuning. The intrinsic history filter shows that history effects extended only 60 ms into the past. The refractory period, corresponding to negative filter coefficients, lasts about 18 ms after a spike. The firing probability then increases and peaks at about 30 ms after a spike.

### 2. Time rescaling for model comparison

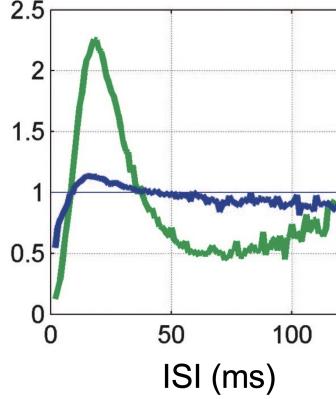To which extent does the inclusion of an intrinsic history filter result in model improvement?



FIG. 24. Mean ratio of observed to expected $z_j$ values indicates that the velocity only model (green) overestimates $\lambda_i(t)$ for about 10 ms after a spike and underestimates $\lambda_i(t)$ for longer times.

For nested models such as considered here, the *time-rescaling theorem* [6, 10] provides a useful tool for addressing this question. Given the spike times $\{t_j\}$ of all spikes emitted by a specific neuron during the $(0, T]$ observational period, we compute for each interspike interval $\tau_j = t_{j+1} - t_j$ the quantity $z_j$ defined as

$$z_j = 1 - \exp\left(-\int_{t_j}^{t_{j+1}} \lambda(t)\, dt\right). \tag{105}$$

For a Poisson process, the interspike interval $\tau_j$ is a random variable with an exponential probability density in the $(0, \infty)$ interval. If $\lambda(t)$ is the correct mean firing rate for this process, $z_j$ will be a random variable with uniform probability density in the $(0, 1)$ interval.

To implement the time-rescaling test, we compute the empirical probability density of the $z_j$ values that result from the observed spike times and the fitted model for $\lambda(t)$, and calculate the ratio of the empirical to the theoretical uniform density. When plotted as a function of the corresponding values of $\tau_j$ (Fig. 24), this ratio clearly demonstrates the modeling

improvement achieved by incorporating the intrinsic spike histoty (blue curve). The velocity model (green curve) overestimates $\lambda_i(t)$ for 10 ms after a spike, and it subsequently underestimates $\lambda_i(t)$. The negative (positive) coefficients of the intrinsic history filter almost completely eliminate the over (under) estimation of the $\lambda_i(t)$ based on the velocity model alone.

## F. Summary

- Generalized linear models provide a principled and systematic approach to modeling the time-dependent rate of inhomogenous Poisson processes that describe the expected firing activity of a neural population.

- The logarithm of the time-dependent rate for each neuron is modeled as a linear combination of intrinsic (the preceding firing activity of all neurons in the population) and extrinsic (the preceding input stimulus or subsequent output activity) observables.

- The likelihood of the data is log-convex; optimal model parameters follow from an unambiguous gradient ascent algorithm.

## III. LECTURE 3: LINEAR AND NONLINEAR DIMENSIONALITY REDUCTION

Let's go back to a motor task that we have already used as an example: the center-out reaching task (Fig. 25). Each trial involves a reach from the center position to one of eight targets on a circle. After target presentation, the subject waits for the go cue to execute to reach. Each trial takes a few seconds.
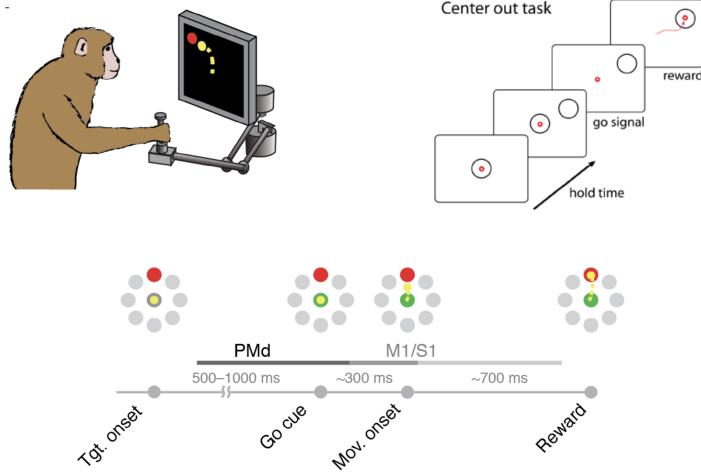


FIG. 25. Instructed delay center-out reaching task

### A. Target dependent population activity

Neural activity in motor and premotor cortices is recorded using multi-electrode arrays (MEAs). Given data for successive reaches to different targets (Fig. 26), we ask whether the pattern of population activity can be decoded to identify the target of each reach. We use a bin of size $\Delta = 300$ ms centered around the time of the go cue to characterize each reach by a point in an empirical neural space of dimension $N$, where $N$ is the number of recorded neurons. Each axis in this neural space represents the firing rate $f_i = y_i/\Delta$ of the $i$th neuron, where $y_i$ is the number of spikes emitted by neuron $i$ during the 300 ms window.
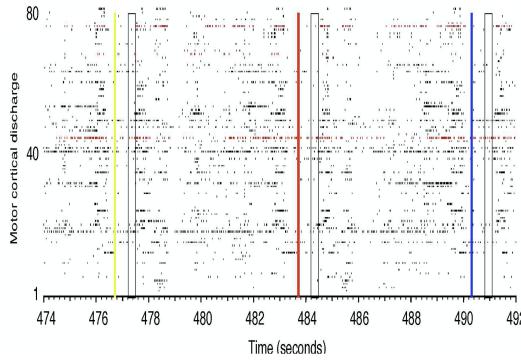


FIG. 26. Activity of $N = 80$ neurons during reaches to three targets. Yellow, red, and blue vertical lines indicate target presentation. Black vertical boxes show 300 ms windows centered around the corresponding go cue.

In order to display the population activity, we choose three neurons (46, 70, and 78; rows shown in red in Fig. 26) and display the resulting three-dimensional neural space. Each point represents a trial, color coded by target (Fig. 27). Given this neural activity, is it possible to infer which target the monkey is going to reach? The answer seems to be NO, as the neural space shows no clustering, no organization according to target!
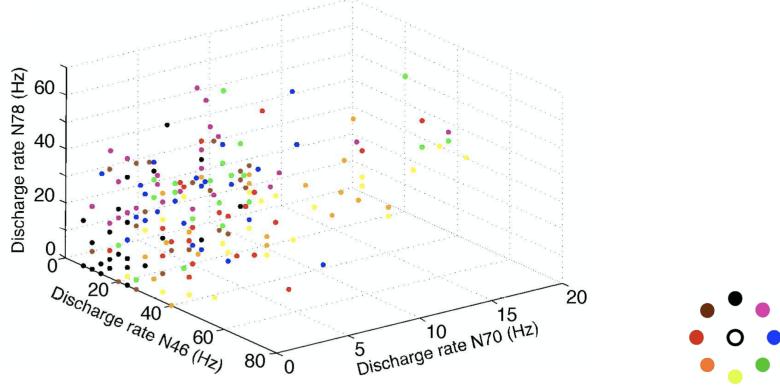
FIG. 27. Activity of the three selected neurons from trial to trial.

This failure suggests that the data should be looked at from a different perspective. Let's recall the simple model of neurons tuned to the direction of motion [8]. In this model, a reach of extent $r$ in the direction $\theta$ would correspond to an activity

$$f_i(r, \theta) = b_i + r\, a_i(1/2)\left[1 + \cos(\theta - \theta_i)\right] , \tag{106}$$

where $b_i$ is the background activity, $a_i$ is the amplitude of activity modulation, and $\theta_i$ is the preferred direction of neuron $i$ (Fig. 21). For a population of $N$ neurons, this model would predict

$$\vec{f}(r, \theta) = \vec{b} + (1/2)\, r\left[\vec{a} + \cos(\theta)\, \vec{\phi}_x + \sin(\theta)\, \vec{\phi}_y\right] , \tag{107}$$

where

$$\begin{aligned}
\vec{b} &= (b_1, b_2, \ldots, b_N) , \\
\vec{a} &= (a_1, a_2, \ldots, a_N) , \\
\vec{\phi}_x &= (a_1 \cos(\theta_1), a_2 \cos(\theta_2), \ldots, a_N \cos(\theta_N)) , \\
\vec{\phi}_y &= (a_1 \sin(\theta_1), a_2 \sin(\theta_2), \ldots, a_N \sin(\theta_N)) .
\end{aligned} \tag{108}$$

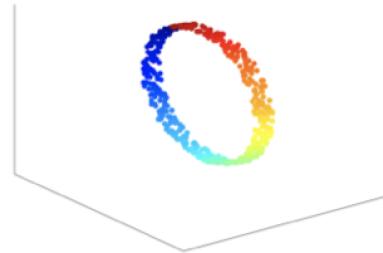A simulation of this model reveals that the data is organized in a low-dimensional ring structure:



FIG. 28. Activity of $N = 100$ orientation selective neurons with preferred orientations uniformly distributed in $[0, 2\pi]$ for planar reaches of extent $10 \leq r \leq 11$.

This observation suggests the use of standard dimensionality reduction methods to find this low-dimensional structure within the $N$-dimensional neural space.

## B. Linear dimensionality reduction: Principal components analysis (PCA)

Consider data in the form of $N$-dimensional vectors.

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{pmatrix} \tag{109}$$

The data here is the $N$-dimensional vector of firing rates associated with each reach. Data for $M$ reaches results in an $N \times M$ matrix

$$X = (\vec{x}_1 | \vec{x}_2 | \dots | \vec{x}_M) = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1M} \\ x_{21} & x_{22} & \cdots & x_{2M} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NM} \end{bmatrix} \tag{110}$$

Estimate the mean firing rate of each neuron:

$$\hat{\mu}_i = \frac{1}{M} \sum_{k=1}^{M} x_{ik} \quad 1 \le i \le N \tag{111}$$

and subtract it from the corresponding row (mean-centered data). Next, use the mean-centered data matrix $X$ to estimate the covariance of the data:

$$\hat{C} = \frac{1}{(M-1)} X X^T \tag{112}$$

$$\hat{C}_{ij} = \frac{1}{(M-1)} \sum_{k=1}^{M} (x_{ik} - \hat{\mu}_i)(x_{jk} - \hat{\mu}_j) \quad \text{for all neurons } 1 \le i,j \le N \tag{113}$$

The diagonalization of the covariance matrix yields eigenvectors and eigenvalues, the principal components:

$$\widehat{C} \vec{u}_\nu = \lambda_\nu \vec{u}_\nu \quad \text{for all } 1 \le \nu \le N \tag{114}$$

### 1. Relation to singular value decomposition

Consider the singular value decomposition (SVD) of the mean-centered data matrix $X$ :

$$X = U \Sigma V^T \tag{115}$$

The columns of the $N \times N$ orthonormal matrix $U$ provide a basis for the neural space. The columns of the $M \times M$ orthonormal matrix $V$ provide a basis for the space of samples. Assume $M > N$. The $N \times M$ matrix $\Sigma$ consists of an $N \times N$ diagonal block and a rectangular block of zeros of size $N \times (M-N)$ to the right of the diagonal block. The matrix $X$ has at most rank $N$; only the leading $N$ columns of $V$, or $N$ rows of $V^T$ contribute to $X$.

Given the singular value decomposition of the data matrix $X$, the eigenvalue decomposition of $XX^T$ is given by

$$\begin{aligned} XX^T &= \left( U\Sigma V^T \right) \left( V \Sigma^T U^T \right) \\ &= U \left( \Sigma \Sigma^T \right) U^T . \end{aligned} \tag{116}$$

The empirical covariance follows:

$$\hat{C} = \frac{1}{(M-1)} X X^T = U \Lambda U^T$$

$$\Lambda = \frac{1}{(M-1)} \Sigma \Sigma^T . \tag{117}$$

## 2. Dimensionality reduction

Consider the diagonal matrix

$$\Lambda = \frac{1}{(M-1)}\Sigma\Sigma^T$$

$$\Lambda = \begin{bmatrix} \lambda_1 & \cdots & 0 & \cdots & 0 \\ 0 & \cdots & \lambda_K & \cdots & 0 \\ 0 & \cdots & 0 & \cdots & \lambda_N \end{bmatrix} \quad \text{with } \lambda_1 \geq \cdots \geq \lambda_K \geq \cdots \geq \lambda_N \ . \tag{118}$$

We keep the $K$ leading eigenvalues, $K < N$ and use the spectral decomposition of the covariance matrix to implement dimensionality reduction (Fig. 29):

$$\hat{C} = \sum_{\mu=1}^{N} \lambda_\mu \vec{u}_\mu \vec{u}_\mu^T \quad \Rightarrow \quad \hat{C} = \sum_{\mu=1}^{K} \lambda_\mu \vec{u}_\mu \vec{u}_\mu^T \tag{119}$$

The space of reduced dimensionality is spanned by the $K$ leading eigenvectors, as described by a matrix $U$ that has the
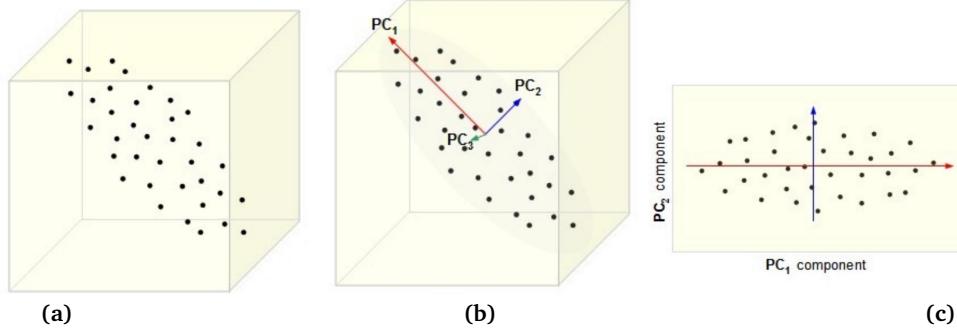


(a)        (b)        (c)

FIG. 29. Dimensionality reduction from a three dimensional original space **(a)** to a two dimensional space spanned by the leading eigenvectors **(b-c)**.

leading eigenvectors $\vec{u}_\nu$, $1 \leq \nu \leq K$ as columns. Then $\hat{C} = U\Lambda U^T$, where $\Lambda$ is the diagonal matrix of $K$ eigenvalues. The coordinates of the data points expressed in the new coordinate system are $Y = U^T X$.

## 3. Principal components as latent variables

We have described how to obtain $Y$ form $X$; similarly, we can reconstruct $X$ from $Y$ as follows (see Fig. 30)

$$X = UY \quad \Rightarrow \quad x_i = \sum_{j=1}^{K} u_{ij} y_j \tag{120}$$

The $i$-th component $u_{ij}$ of the $j$-th eigenvector is the "weight" from $y_j$ to $x_i$. We can interpret this algorithm as a generative model in which the full-dimension data matrix $X$ is constructed from the low-dimensional data matrix $Y$.

$$P(y_j) = \mathcal{N}(0, \lambda_j) \quad \rightarrow \quad x_i = \sum_{j=1}^{d} u_{ij} y_j \tag{121}$$

Two extensions of PCA allow for the inclusion of noise in the generative model. Probabilistic PCA (PPCA) introduces
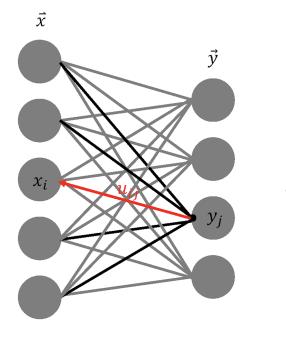
FIG. 30. Principal components as latent variables in a generative model.

homogeneous noise into the generative process:

$$x_i = \sum_{j=1}^{d} u_{ij} y_j + \eta_i$$

$$P(\eta_i) = \mathcal{N}\left(0, \sigma^2\right)$$

$$P(y_j) = \mathcal{N}\left(0, \lambda_j\right) \tag{122}$$

$$P(x_i \mid y_j) = \mathcal{N}\left(u_{ij} y_j, \sigma^2\right)$$

$$P(x_i \mid \vec{y}) = \mathcal{N}\left(\sum_j u_{ij} y_j, \sigma^2\right)$$



FIG. 31. Neural activity in PMd preceding reaches in a center-out task. The three-dimensional representation is a projection onto a subspace spanned by the three leading principal components (neural modes).

The constraint of homogeneous noise can be relaxed to allow for inhomogeneous noise whose amplitude depends on $i$. This leads from PPCA to Factors Analysis (FA):

$$P(\eta_i) = \mathcal{N}\left(0, \sigma_i^2\right)$$

$$P(x_i \mid y_j) = \mathcal{N}\left(u_{ij} y_j, \sigma_i^2\right) \tag{123}$$

$$P(x_i \mid \vec{y}) = \mathcal{N}\left(\sum_j u_{ij} y_j, \sigma_i^2\right)$$

Factor Analysis has been applied to neural data recorded from dorsal premotor cortex (PMd) during the execution of a center-out task [11]. This is a region where neural activity is known to correlate with the endpoint of an upcoming reach,

including both direction and extent. In contrast to the random three-dimensional representation of Fig. 27, the projection of the data onto a three-dimensional space spanned by the three leading principal components reveals a structure that reflects the separation and spatial organization of the targets, see Fig. 31.

## C. Dimensionality reduction: Neural modes and latent variables

Let's consider a schematic scenario for the recording of neural population activity using multi-electrode arrays (Fig. 32). The recording device subsamples the population activity and defines an empirical neural space of dimension $N$, the number of recorded neurons. Each axis in the empirical neural space corresponds to the firing rate of one of the recorded neurons. At every time bin of size $\Delta$, the population activity is described by a point in this empirical neural space. As time evolves, consecutive points define a trajectory that has been consistently found to be confined to a low-dimensional manifold within the empirical neural space [12]. The observed low-dimensionality of population dynamics is not only restricted to the motor
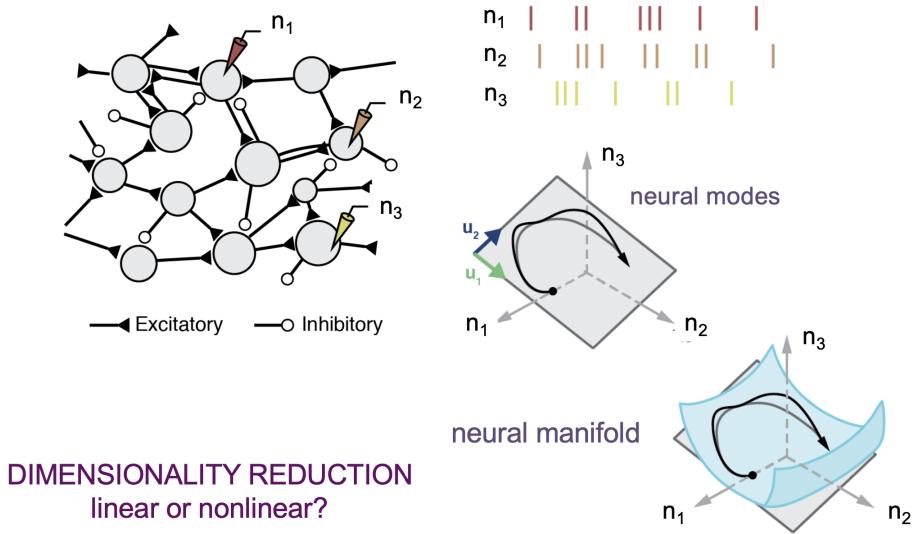


FIG. 32. Subsampling of the population neural activity defines an empirical neural space. The trajectory that describes the time evolution of the neural activity is confined to a low-dimensional manifold that can be either linear or nonlinear.

cortex, it is also observed in many other brain areas: frontal, prefrontal, parietal, visual, auditory, and olfactory cortices (see [13] and references therein). A two-dimensional neural manifold that encodes for both position and accumulation of evidence has recently been characterized in hippocampus [14].

An important aspect of these neural manifolds is whether they are linear or nonlinear [15], an issue that leads to the concepts of *intrinsic dimension* and *embedding dimension* [16]. As illustrated in Fig. 33, the intrinsic dimension is the true dimension of a nonlinear manifold, while the embedding dimension is the dimensionality of the smallest flat space that fully contains the manifold. The difference between these two dimensions is a proxy for the degree of manifold nonlinearity [15]. A useful approach to manifold identification is to use PCA to obtain a flat space whose dimensionality is low but not too low; this dimensionality must be large enough for the flat space to play the role of an embedding space. More sphisticated methods for nonlinear dimensionality reduction can then be used to identify the nonlinear manifold within the flat embedding space.

## D. Nonlinear dimensionality reduction: Isomap and multidimensional scaling

Let's go back to the center-out task that we introduced at the beginning of this lecture. If we apply PCA to data such as shown in Fig. 26, we find eigenvalues that gradually decrease as the dimensionality $K$ of the embedding space increases (Fig. 34, red curve). This gradual decay does not obviously define a threshold that identifies a specific value for $K$. In contrast, if we apply a nonlinear method for dimensionality reduction such as Isomap [17] to the same data (Fig. 34, blue curve) we find that the first two eigenvalues dominate the others, which stay at an almost constant level. The magnitude
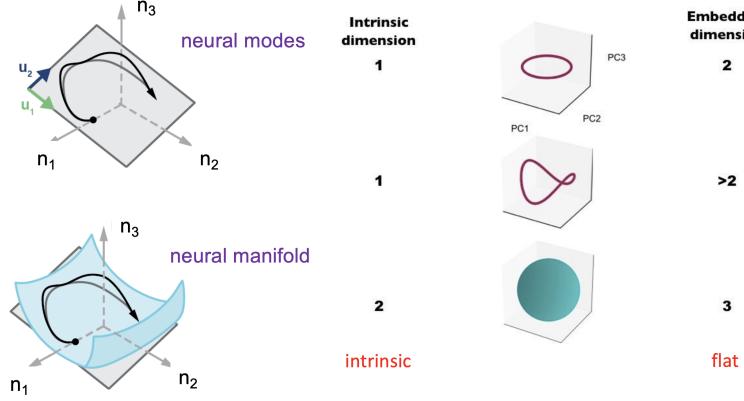
FIG. 33. Intrinsic vs embedding dimension.

of the eigenvalues $\lambda_i$ for $i \geq 3$ corresponds to the amplitude of the noise fluctuations that give the data cloud thickness in all nonsignificant dimensions. The leading dimensions describe a nonlinear manifold of intrinsic dimension two.
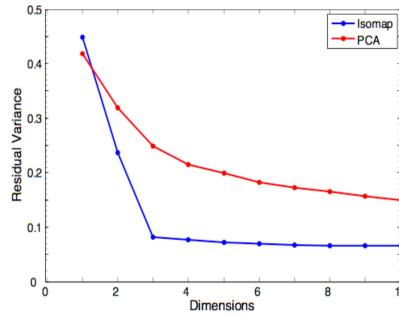


FIG. 34. Eigenvalue spectrum for center-out task data obtained by PCA (red) and Isomap (blue).

How does Isomap work [17]? It flattens the curved manifold into a flat manifold of the same dimension (Fig. 35). Isomap is based on MultiDimensional Scaling (MDS), a method devised to find a solution to the following problem (Fig.
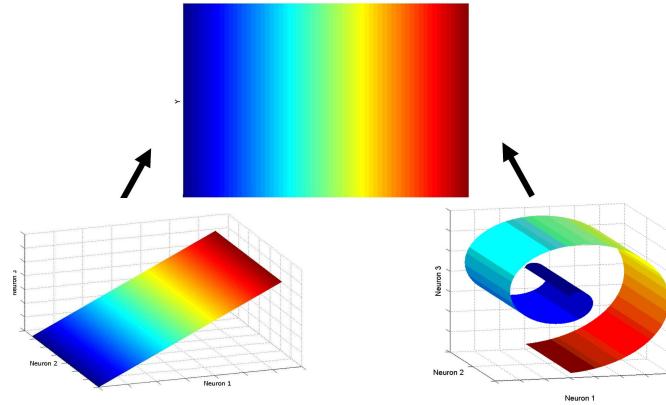


FIG. 35. Isomap flattens curved manifolds.

36): given an $M \times M$ matrix of dissimilarities between $M$ objects, is it possible to represent these objects as $M$ points in an Euclidean space such that the matrix of Euclidean distances between them matches as well as possible the original dissimilarity matrix [18]?

Consider data in the form of $N$-dimensional vectors, such as an $N$-dimensional vector of firing rates associated with each

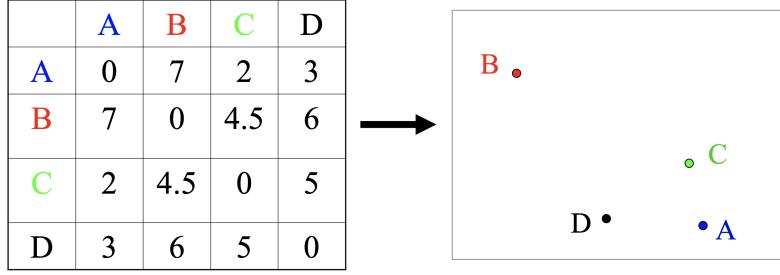|   | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 7 | 2 | 3 |
| B | 7 | 0 | 4.5 | 6 |
| C | 2 | 4.5 | 0 | 5 |
| D | 3 | 6 | 5 | 0 |

FIG. 36. Matrix of dissimilarities between $M = 4$ objects (left); representation of these objects as points in an Euclidean space (right).

reach. Data for $M$ reaches result in an $N \times M$ data matrix $X$:

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{pmatrix} \implies X = (\vec{x}_1 | \vec{x}_2 | \dots | \vec{x}_M) \tag{124}$$

If the matrix $X$ is hidden from us, but we are given instead an $M \times M$ matrix $S$ of squared distances between the points, can we reconstruct the matrix $X$? Let's start by considering the case where the original distances are Euclidean:

$$S_{ij} = d_{ij}^2 = \left(\vec{x}_i - \vec{x}_j\right)^T \left(\vec{x}_i - \vec{x}_j\right) . \tag{125}$$

The scalar product between data points can be written as:

$$\vec{x}_i^T \vec{x}_j = -(1/2)\left(S_{ij} - \|\vec{x}_i\|^2 - \|\vec{x}_j\|^2\right) . \tag{126}$$

In matrix form,

$$\vec{x}_i^T \vec{x}_j = \left(X^T X\right)_{ij} , \tag{127}$$

and

$$S_{ij} - \|\vec{x}_i\|^2 - \|\vec{x}_j\|^2 = (JSJ)_{ij} , \tag{128}$$

where $J$ is the $M \times M$ centering matrix $J = I - (1/M)\, ee^T$ with $e$ being a column vector with all entries equal to 1. Then

$$X^T X = -(1/2)\, JSJ . \tag{129}$$

From this equation the data matrix $X$ can be easily obtained:

$$X^T X = U \Lambda U^T \implies X = \Lambda^{1/2} U^T . \tag{130}$$

If the distance matrix to which this calculation is applied is based on Euclidean distances, this process allows us to recover the data matrix $X$ from the matrix $S$ of squared distances. A reduction of the dimensionality of the Euclidean representation follows from truncating of the number of eigenvalues in the diagonal matrix $\Lambda$ from $M$ to $K$; this restricts to $K$ the number of eigenvectors used to reconstruct $X$. It can be proved that this truncation is equivalent to PCA, which is based on the diagonalization of $XX^T$.

When applied to an arbitrary matrix $S$ of squared distances, the method follows similar steps. An 'inner product' matrix $Y$ is defined through the centering operation: $Y = -(1/2)\, JSJ$, followed by the diagonalization of $Y$: $Y = U \Lambda U^T$ and the identification of the data matrix $X$ as $X = \Lambda^{1/2} U^T$. This procedure minimizes a cost function $E$ that measures the Frobenius norm of the difference between two matrices: the matrix $Y$ defined by the original $S$ and the inner product matrix $X^T X$ obtained from the Euclidean representation of the data:

$$E(X) = \left\|X^T X - Y\right\|_F . \tag{131}$$

The beauty of the Isomap approach is based on the proposed method for computing the matrix $S$ of square distances, which is obtained from the data by constructing empirical Geodesics [17]. For each data point, a number $P$ of nearest neighbors is chosen, and the short distances from the chosen point to its $P$ neighbors are computed as Euclidean. The
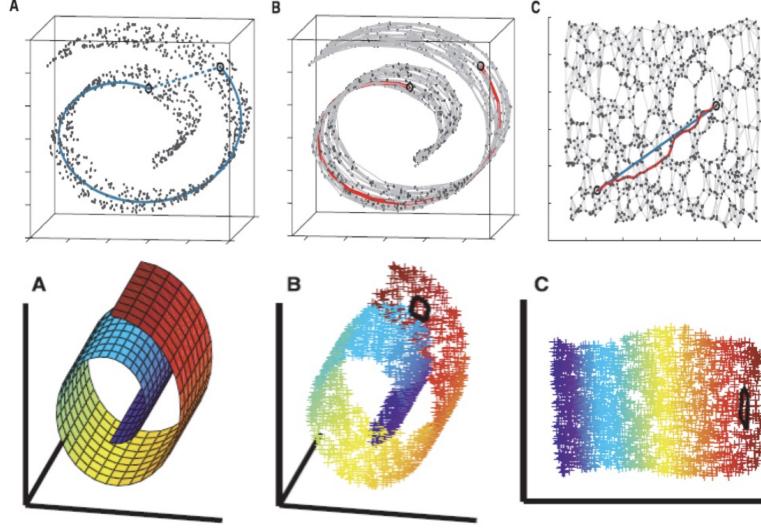
FIG. 37. The distance between two points that are not nearest neighbors is defined through an empirical Geodesic: the shortest path that connects the two points through data points. This Geodesic distance will typically be larger than the Euclidean distance between the two points (see dashed blue line in A).

distance between two points that are not nearest neighbors follows from walking in the manifold. Every step in this path involves a known short Euclidean distance between neighboring points; the shortest path that connects the two end points defines their Geodesic distance (Fig. 37, top row) .

When applied to the Swiss roll data [17], Isomap results in a flat two-dimensional representation that preserves the neighborhood and distance relations among data points (Fig. 37, bottom row). When applied to neural data recorded during the execution of the center-out task, the Isomap eigenvalue spectrum (Fig. 34, blue curve) signals a two-dimensional curved manifold. The Isomap projection then provides a flat two-dimensional clustered representation that captures the spatial organization of the targets (Fig. 38). The degree of nonlinearity of the two-dimensional manifold can be inferred
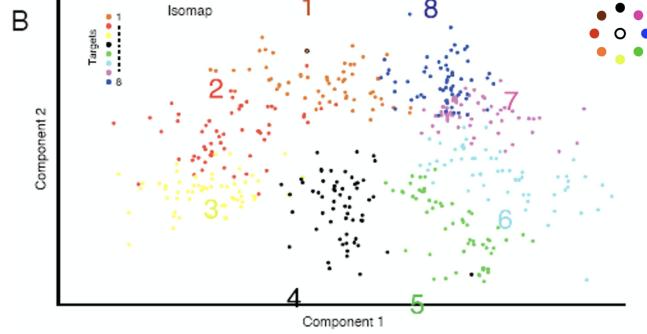


FIG. 38. Two dimensional representation of neural data for the center-out task shows clusters that capture the spatial organization of the targets.

from a comparison of pairwise Geodesic and Euclidean distances (Fig. 39). Note that the distinction between these two distances vanishes for the short distances between nearest neighbors, and increases monotonically for larger distances.

### E.   Summary

- The dynamics of neural population activity is typically confined to low-dimensional manifolds within a high-dimensional neural space.

- These manifolds can be linear or nonlinear. A nonlinear manifold is characterized by its *intrinsic dimension* and by the
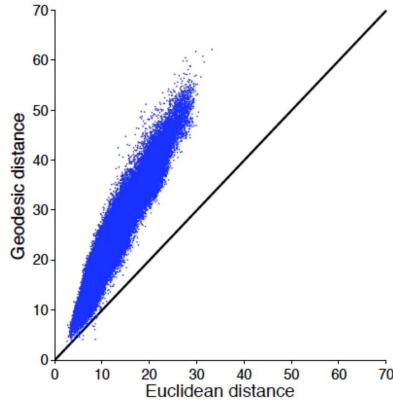
FIG. 39. Geodesic vs Euclidean distances.

*embedding dimension* of the smallest flat space that fully contains it. The difference between these two dimensions is a proxy for the degree of manifold nonlinearity.

- MultiDimensional Scaling in combination with the empirical Geodesic distances proposed by Isomap provide a useful tool for obtaing flat representations of nonlinear manifolds.

### ACKNOWLEDGMENTS

[1] B. Widrow and R. Winter, Neural nets for adaptive filtering and adaptive pattern recognition, IEEE Computer Magazine **21**, 25 (1988).
[2] E. Levin, N. Tishby, and S. A. Solla, A statistical approach to learning and generalization in neural networks, Proceedings of the IEEE **78**, 1568 (1990).
[3] S. A. Solla, A Bayesian approach to learning in neural networks, International Journal of Neural Systems **6**, 161 (1995).
[4] S. A. Solla and E. Levin, Learning in neural networks: The validity of the annealed approximation, Physical Review A **46**, 2124 (1992).
[5] P. McCullagh and J. A. Nelder, *Generalized linear models* (Chapman & Hall, 2019).
[6] W. Truccolo, U. Eden, M. Fellows, J. Donoghue, and E. Brown, A point process framework for relating neural spiking activity to spiking history, neural ensemble, and extrinsic covariate effects, Journal of Neurophysiology **93**, 1074 (2005).
[7] W. Truccolo, L. R. Hochberg, and J. P. Donoghue, Collective dynamics in human and monkey sensorimotor cortex: predicting single neuron spikes, Nature Neuroscience **13**, 105 (2010).
[8] A. P. Georgopoulos, J. F. Kalaska, R. Caminiti, and J. T. Massey, On the relations between the direction of two-dimensional arm movements and cell discharge in primate motor cortex, Journal of Neuroscience **2**, 1527 (1982).
[9] L. Paninski, M. R. Fellows, N. G. Hatsopoulos, and J. P. Donoghue, Spatiotemporal tuning of motor cortical neurons for hand position and velocity, Journal of Neurophysiology **91**, 515 (2004).
[10] E. Brown, R. Barbieri, V. Ventura, R. Kass, and L. Frank, The time-rescaling theorem and its application to neural spike train data analysis, Neural Computation **14**, 325 (2002).
[11] G. Santhanam, B. M. Yu, V. Gilja, S. I. Ryu, A. Afshar, M. Sahani, and K. V. Shenoy, Factor-Analysis methods for higher-performance neural prostheses, Journal of Neurophysiology **102**, 1315 (2009).
[12] J. P. Cunningham and B. M. Yu, Dimensionality reduction for large-scale neural recordings, Nature Neuroscience **17**, 1500 (2014).
[13] J. A. Gallego, M. G. Perich, L. E. Miller, and S. A. Solla, Neural manifolds for the control of movement, Neuron **94**, 978 (2017).
[14] E. H. Nieh, M. Schottdorf, N. W. Freeman, R. J. Low, S. Lewallen, S. A. Koay, L. Pinto, J. L. Gauthier, C. D. Brody, and D. W. Tank, Geometry of abstract learned knowledge in the hippocampus, Nature **595**, 80 (2021).
[15] E. Altan, S. A. Solla, L. E. Miller, and E. J. Perreault, Estimating the dimensionality of the manifold underlying multi-electrode neural recordings, PLoS Computational Biology **17**, e1008591 (2021).
[16] M. Jazayeri and S. Ostojic, Interpreting neural computations by examining intrinsic and embedding dimensionality of neural activity, Current Opinion in Neurobiology **70**, 113 (2021).

[17] J. B. Tenenbaum, V. de Silva, and J. C. Langford, A global geometric framework for nonlinear dimensionality reduction, Science **290**, 2319 (2000).

[18] T. F. Cox and M. A. A. Cox, *Multidimensional scaling* (Chapman & Hall, 2000).