


# Learning Quantitative Sequence–Function Relationships from Massively Parallel Experiments

Gurinder S. Atwal<sup>1</sup> · Justin B. Kinney<sup>1</sup> 

Received: 29 May 2015 / Accepted: 13 October 2015 / Published online: 29 October 2015  
© The Author(s) 2015. This article is published with open access at Springerlink.com

**Abstract** A fundamental aspect of biological information processing is the ubiquity of sequence–function relationships—functions that map the sequence of DNA, RNA, or protein to a biochemically relevant activity. Most sequence–function relationships in biology are quantitative, but only recently have experimental techniques for effectively measuring these relationships been developed. The advent of such “massively parallel” experiments presents an exciting opportunity for the concepts and methods of statistical physics to inform the study of biological systems. After reviewing these recent experimental advances, we focus on the problem of how to infer parametric models of sequence–function relationships from the data produced by these experiments. Specifically, we retrace and extend recent theoretical work showing that inference based on mutual information, not the standard likelihood-based approach, is often necessary for accurately learning the parameters of these models. Closely connected with this result is the emergence of “diffeomorphic modes”—directions in parameter space that are far less constrained by data than likelihood-based inference would suggest. Analogous to Goldstone modes in physics, diffeomorphic modes arise from an arbitrarily broken symmetry of the inference problem. An analytically tractable model of a massively parallel experiment is then described, providing an explicit demonstration of these fundamental aspects of statistical inference. This paper concludes with an outlook on the theoretical and computational challenges currently facing studies of quantitative sequence–function relationships.

**Keywords** Sequence–function relationships · Mutual information · Likelihood · Diffeomorphic modes · Sort-Seq

---

✉ Justin B. Kinney  
jkinney@cshl.edu

<sup>1</sup> Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA

## 1 Introduction

A major long-term goal in biology is to understand how biological function is encoded within the sequences of DNA, RNA, and protein. The canonical success story in this effort is the genetic code: given an arbitrary sequence of messenger RNA, the genetic code allows us to predict with near certainty what peptide sequence will result. There are many other biological codes we would like to learn as well. How does the DNA sequence of a promoter or enhancer encode transcriptional regulatory programs? How does the sequence of pre-mRNA govern which exons are kept and which are removed from the final spliced mRNA? How does the peptide sequence of an antibody govern how strongly it binds to target antigens?

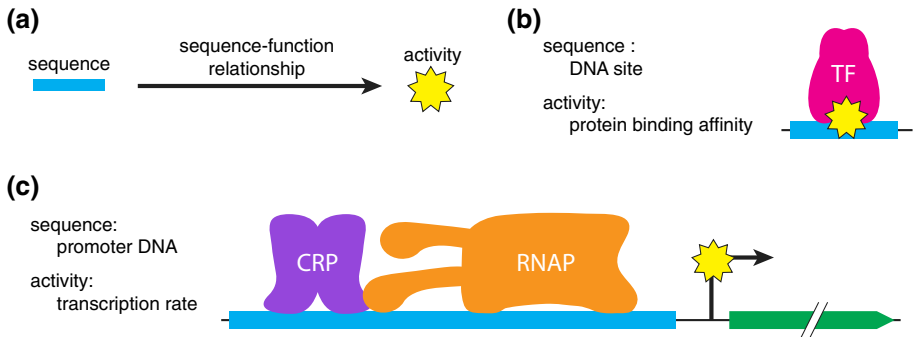
A major difference between the genetic code and these other codes is that while the former is qualitative in nature, the latter are governed by sequence–function relationships that are inherently quantitative. Quantitative sequence–function relationships<sup>1</sup> describe any function that maps the sequence of a biological heteropolymer to a biologically relevant activity (Fig. 1a). Perhaps the simplest example of such a relationship is how the affinity of a transcription factor protein for its DNA binding site depends on the DNA sequence of that site (Fig. 1b). Such relationships are a key component of the more complicated relationship between the DNA sequence of a promoter or enhancer (which typically binds multiple proteins) and the resulting rate of mRNA transcription (Fig. 1c). In both of these cases, the activities of interest (affinity or transcription rate) can vary over orders of magnitude and yet still be finely tuned by adjusting the corresponding sequence (binding site or promoter/enhancer). Similarly, other sequence–function relationships, like the inclusion of exons during mRNA splicing or the affinity of a protein for its ligand, are fundamentally quantitative.

The study of quantitative sequence–function relationships presents an exciting opportunity for the concepts and methods of statistical physics to shed light on biological systems. There is a natural analogy between biological sequences and the microstates of physical systems, as well as between biological activities and physical Hamiltonians. Yet we currently lack answers to basic questions a statistical physicist might ask, such as “what is the density of states?” or “is a relationship convex or glassy?” The answers to such questions may well have important consequences for diverse fields including biochemistry, systems biology, immunology, and evolution.

Experimental methods for measuring sequence–function relationships have improved dramatically in recent years. In the mid 2000s, multiple “high-throughput” methods for measuring the DNA sequence specificity of transcription factors were developed; these methods include protein binding microarrays (PBMs) [2,3], *Escherichia coli* one-hybrid technology (E1H) [4], and microfluidic platforms [5]. The subsequent development and dissemination of ultra-high-throughput DNA sequencing technologies then led, starting in 2009, to the creation of a number of “massively parallel” experimental techniques for probing a wide range of sequence–function relationships (Table 1). These massively parallel assays can readily measure the functional activity of  $10^3$  to  $10^8$  sequences in a single experiment by coupling standard bench-top techniques to ultra-high-throughput DNA sequencing.

Massively parallel experiments are very unlike conventional experiments in physics: they are typically very noisy and rarely provide direct readouts of the quantities that one cares about. Moreover, the noise characteristics of these measurements are difficult to accurately model. Indeed, such noise generally exhibits substantial day-to-day variability. Although standard inference methods require an explicit model of experimental noise, it is still possible

<sup>1</sup> These have also been called quantitative sequence-activity maps, or QSAMs [1].



**Fig. 1** Sequence–function relationships in biology. **a** A sequence–function relationship maps a biological sequence (blue bar) to a biologically relevant activity (yellow star). **b** One of the simplest sequence–function relationships is how the affinity (star) of a transcription factor protein (magenta) for its DNA binding site depends on the sequence of that site (blue). **c** A more complicated sequence–function relationship describes how the rate of mRNA transcription depends on the DNA sequence of a gene’s promoter region. At the *lac* promoter of *E. coli*, this transcription rate (star) depends on how strongly both the transcription factor CRP (purple) and the RNA polymerase holoenzyme (RNAP; orange) bind their respective sites within the promoter region (blue)

to precisely learn quantitative sequence–function relationships from massively parallel data even when noise characteristics are unknown [27,28].

The ability to fit parametric models to these data reflects subtle but important distinctions between two objective functions used for statistical inference: (i) likelihood, which requires *a priori* knowledge of the experimental noise function and (ii) mutual information [29], a quantity based on the concept of entropy, which does not require a noise function. In contrast to the conventional wisdom that more experimental measurements will improve the model inference task, the standard maximum likelihood approach will *typically* never learn the right model, even in the infinite data limit, if one uses an imperfect model of experimental noise. Model inference based on mutual information does not suffer from this ailment.

Mutual-information-based inference is unable to pin down the values of model parameters along certain directions in parameter space known as “diffeomorphic modes” [28]. This inability is not a shortcoming of mutual information, but rather reflects a fundamental distinction between how diffeomorphic and nondiffeomorphic directions in parameter space are constrained by data. Analogous to the emergence of Goldstone modes in particle physics due to a specific yet arbitrary choice of phase, diffeomorphic modes arise from a somewhat arbitrary choice of the sequence-dependent activity that one wishes to model. Likelihood, in contrast to mutual information, is oblivious to the distinction between diffeomorphic and nondiffeomorphic modes.

We begin this paper by briefly reviewing a variety of massively parallel assays for probing quantitative sequence–function relationships. We then turn to the problem of learning parametric models of these relationships from the data that these experiments generate. After reviewing recent work on this problem [28], we extend this work in three ways. First, we show that “diffeomorphic modes” of the parametric activity model that one wishes to learn are “dual” to certain transformations of the corresponding model of experimental noise (the “noise function”). This duality reveals a symmetry of the inference problem, thereby establishing a close analogy with Goldstone modes. Next we compute and compare the Hessians of likelihood and mutual information. This comparison suggests an additional analogy between this inference problem and concepts in fluid mechanics. Finally, we work through an analytically tractable model of a massively parallel experiment of protein–DNA binding. This

**Table 1** Massively parallel experiments used for studying various sequence–function relationships

Sequence	Activity	System	Name	Publication	
DNA binding sites	Protein–DNA binding affinity	Purified protein	Bind-n-Seq	Zykovich et al. [6]	
			HT-SELEX	Zhao et al. [7]	
			EMSA-Seq	Jolma et al. [8]	
			SELEX-Seq	Wong et al. [9]	
Promoter/ enhancer DNA	Transcription rate	Purified protein		Patwardhan et al. [11]	
			Bacteria	Sort-Seq	Kinney et al. [12]
			Cell culture	MPRA	Melnikov et al. [1]
			Mouse liver		Patwardhan et al. [13]
			Yeast		Sharon et al. [14]
		Mouse retina	CRE-Seq	Kwasniesk et al. [15]	
Protein	Ligand binding	Phage display	DMS	Fowler et al. [16]	
	Cellular growth rate	Yeast	EMPIRIC	Hietpas et al. [17]	
	Toxin activity	Bacteria		Adkar et al. [18]	
	H1N1 binding	Yeast display		Whitehead et al. [19]	
	GPCR expression	Bacteria		Schlinkmann et al. [20]	
RNA	mRNA translation	Bacteria		Holmqvist et al. [21]	
	sRNA targeting	Bacteria	qSortSeq	Peterman et al. [22]	
	mRNA translation	Cell culture		Oikonomou et al. [23]	
	mRNA translation	Cell culture	FACS-Seq	Noderer et al. [24]	
Replication origins	DNA replication	Yeast	ARS-Seq	Liachko et al. [25]	
Endonuclease sites	DNA cutting	Purified protein		Thyme et al. [26]	

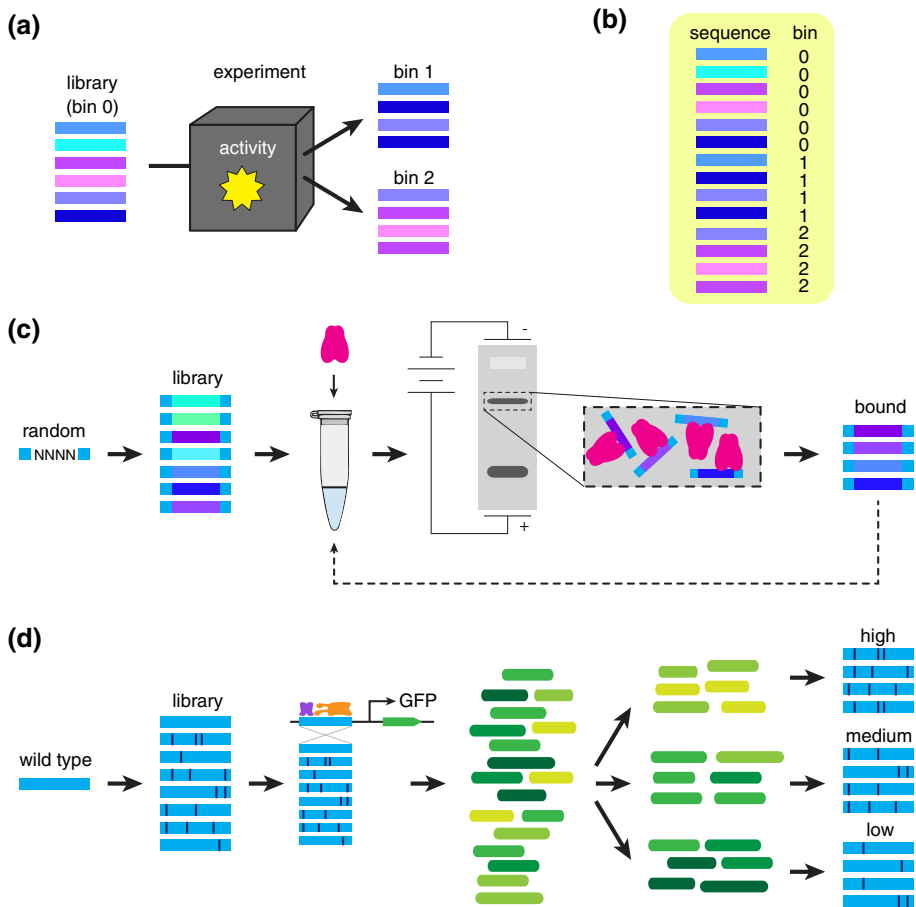
Columns show the type of sequences interrogated, the sequence activity assayed, the biological system on which the experiments were performed, the name (if any) of the experimental technique, and the publication describing the method. This table is not exhaustive; it only describes some of the earliest experiments in each type of system

example explicitly illustrates the differences between likelihood- and mutual-information-based approaches to inference, as well as the emergence of diffeomorphic modes.

It should be noted that the inference of receptive fields in sensory neuroscience is another area of biology in which mutual information has proved useful as an objective function, and that work in this area has also provided important insights into basic aspects of machine learning [30–34]. Indeed, the problem of learning quantitative sequence–function relationships in molecular biology is very similar to the problem of learning receptive fields in neuroscience [28]. The discussion of this problem in the neuroscience context, however, has largely avoided in-depth analyses of how mutual information relates to likelihood, as well as of how diffeomorphic modes emerge.

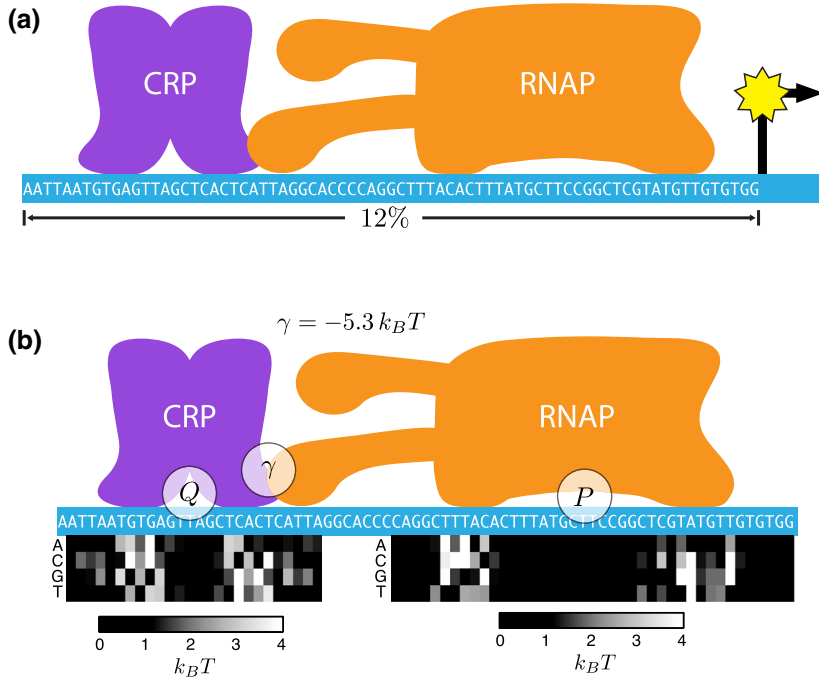
## 2 Massively Parallel Experiments Probing Sequence–Function Relationships

All of the massively parallel experiments in Table 1 share a common structure (Fig. 2a). The first step in each experiment is to generate a large set of (roughly  $10^3$  to  $10^8$ ) different



**Fig. 2** Overview of massively parallel experiments for studying quantitative sequence–function relationships. **a** The input to each experiment is a library of different sequences that one wishes to test. The output is one or more bins of sequences; each sequence in each bin is randomly selected from the library with a weight that depends on a measurement of that sequence’s activity (*star*). **b** The resulting data set consists of a list of (non-unique) sequences, each sequence assigned to either the input library or one of the output bins. **c** Illustration of experimental methods for measuring the sequence-dependent binding energy of purified transcription factor proteins. The input library typically consists of random DNA flanked by constant sequence. This library DNA is mixed with the protein of interest and binding is allowed to come to equilibrium. DNA bound by protein is then separated from unbound DNA, e.g. by running complexes on a gel (shown), then sequenced along with a sample from the input library. **d** Sort-Seq [12] is a massively parallel experiment that uses a library of mutagenized sequences to probe the mechanisms of transcriptional regulation employed by a specific wild type promoter of interest. Mutant promoters are cloned upstream of the GFP gene, and *E. coli* cells harboring these expression constructs are sorted into bins using FACS. The mutant promoters in each bin, as well as promoters from the input library, are then sequenced

sequences to measure. This set of sequences is called the “library.” Multiple different types of libraries can be used depending on the application. One then performs an experiment that takes this library as input, and as output provides a set of one or more “bins” of sequences. Each output bin contains sequences selected from the library with a weight that depends on the measured activity of that sequence. Finally, a sample of sequences from each of the



**Fig. 3** The *lac* promoter region studied in [12]. **a** Sort-Seq was used to dissect a 75 bp region of the *E. coli lac* promoter using a library consisting of wild type sequences mutagenized at 12% per nucleotide, i.e., each library sequence had nine mutations on average. **b** The resulting data were used to learn a quantitative sequence–function relationship, the mathematical form of which reflected an explicit biophysical model of transcriptional regulation. This model included two “energy matrices” describing the sequence-dependent binding energy of CRP ( $Q$ ) and RNAP ( $P$ ) to their respective sites. It also included a value for the interaction energy  $\gamma$  between these two proteins

output bins, as well as from the input library, are determined using ultra-high-throughput DNA sequencing. The resulting data thus consists of a long list of (typically non-unique) DNA sequences, each assigned to a corresponding bin (Fig. 2b). It is from these data that we wish to learn quantitative models of sequence–function relationships.

Some of the earliest massively parallel experiments were designed to measure the specificity of purified transcription factors for their DNA binding sites [6–10] (Fig. 2c). The library used in such studies consists of a fixed-length region of random DNA flanked by constant sequences used for PCR amplification. This library is mixed with the transcription factor of interest, after which protein–bound DNA is separated from unbound DNA, e.g., by running the protein–DNA mixture on a gel. Protein–bound DNA is then sequenced, along with the input library.

Using a library of random DNA to assay protein–DNA binding has the advantage that the same library can be used to study each protein. This is particularly useful when performing assays on many different proteins at once (e.g., [8, 35]). On the other hand, only a very small fraction of library sequences will be specifically bound by the protein of interest. Moreover,

because proteins typically bind DNA in a non-specific manner, such experiments are often performed serially in order to achieve substantial enrichment.<sup>2</sup>

The first massively parallel experiment to probe how multi-protein–DNA complexes regulate transcription in living cells was Sort-Seq [12] (Fig. 2d). The sequence library used in this experiment was generated by introducing randomly scattered mutations into a “wild type” sequence of interest, specifically, the 75 bp region of the promoter of the *lac* gene in *E. coli* depicted in Fig. 3a. A few million of these mutant promoters were cloned upstream of the green fluorescent protein (GFP) gene. Cells carrying these expression constructs were grown under conditions favorable to promoter activity and were then sorted into a small number of bins according to each cell’s measured fluorescence. This partitioning of cells was accomplished using fluorescence-activated cell sorting (FACS) [41], a method that can readily sort  $\sim 10^4$  cells per second. The mutant promoters within each sorted bin as well as within the input library were then sequenced, yielding measurements for  $\sim 2 \times 10^5$  variant promoter sequences. We note that advances in DNA sequencing have since made it possible to accumulate much more data, and it is no longer difficult to measure the activities of  $\sim 10^7$  different sequences in this manner.

Massively parallel experiments using mutagenized sequences provide data about sequence–function relationships within a localized region of sequence space centered on the wild type sequence of interest. Measuring these local relationships can provide a wealth of information about the functional mechanisms of the wild type sequence. For instance, the Sort-Seq data of [12] allowed the inference of an explicit biophysical model for how CRP and RNAP work together to regulate transcription at the *lac* promoter (Fig. 3b). In particular, the authors used their data to learn quantitative models for the in vivo sequence specificity of both CRP and RNAP. Model fitting also enabled measurement of the protein–protein interaction by which CRP is able to recruit RNAP and up-regulate transcription.

Mutagenized sequences have also been used extensively for “deep mutational scanning” experiments on proteins. In this context, selection experiments on mutagenized proteins allow one to identify protein domains critical for folding and function. A variety of deep mutational scanning experiments are described in [42].

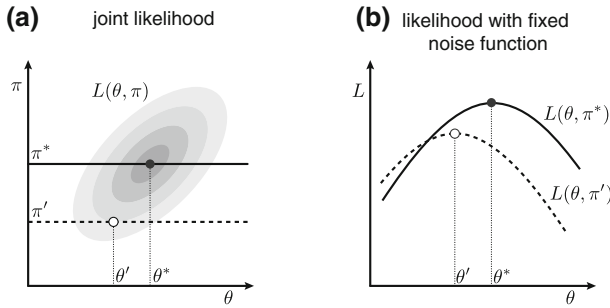
### 3 Inference Using likelihood

The inference of quantitative sequence–function relationships from massively parallel experiments can be phrased as follows. Data consists of a large number of sequences  $\{S_n\}_{n=1}^N$ , each sequence  $S$  having a corresponding measurement  $M$ . Due to experimental noise, repeated measurements of the same sequence  $S$  can yield different values for  $M$ . Our experiment therefore has the following probabilistic form:

$$\underset{\text{sequence}}{S} \xrightarrow[\text{experiment}]{p(M|S)} \underset{\text{measurement}}{M} \quad (1)$$

If we assume that the measurements for each sequence are independent, and if we have an explicit parametric form for  $p(M|S)$ , then we can learn the values of the parameters by maximizing the per-datum log likelihood,

<sup>2</sup> This serial enrichment approach is known as SELEX and is much older than ultra-high-throughput DNA sequencing; see [36–40].



**Fig. 4** Schematic illustration of how likelihood  $L(\theta, \pi)$  depends on the model  $\theta$  and the noise function  $\pi$  in the  $N \rightarrow \infty$  limit. **a, b**  $L$  will typically have a correlated dependence on  $\theta$  and  $\pi$ . If  $\pi$  is set equal to the correct noise function  $\pi^*$ , then  $L$  will be maximized by the correct model  $\theta^*$ . However, if  $\pi$  is set to an incorrect noise function  $\pi'$ ,  $L$  will typically attain a maximum at an incorrect  $\theta'$

$$L = \frac{1}{N} \sum_{n=1}^N \log p(M_n | S_n). \tag{2}$$

In what follows we will refer to the quantity  $L$  simply as the “likelihood.”

In regression problems such as this, one introduces an additional layer of structure. Specifically, we assume the measurement  $M$  of each sequence  $S$  is a noisy readout of some underlying activity  $R$  that is a deterministic function of that sequence. We call the function relating  $R$  to  $S$  the “activity model” and denote it using  $\theta(S)$ . This activity model is ultimately what we want to understand. The specific way the activity  $R$  is read out by measurements  $M$  is then specified by a conditional probability distribution,  $\pi(M|R)$ , which we call the “noise function.”<sup>3</sup> Our experiment is thus represented by the Markov chain

$$\begin{array}{ccccc}
 S & \xrightarrow{\text{model}} & R & \xrightarrow{\text{noise function}} & M \\
 \text{sequence} & \theta(S) & \text{activity} & \pi(M|R) & \text{measurement}
 \end{array} \tag{3}$$

The corresponding likelihood is

$$L(\theta, \pi) = \frac{1}{N} \sum_{n=1}^N \log \pi(M_n | \theta(S_n)). \tag{4}$$

The model we adopt for our experiment therefore has two components:  $\theta$ , which describes the sequence–function relationship of interest, and  $\pi$ , which we do not really care about.

Standard statistical regression requires that the noise function  $\pi$  be specified up-front.  $\pi$  can be learned either by performing separate calibration experiments, or by assuming a functional form based on an educated guess. This can be problematic, however. Consider inference in the large data limit,  $N \rightarrow \infty$ , which is illustrated in Fig. 4. Likelihood is determined by both the model  $\theta$  and the noise function  $\pi$  (Fig. 4a). If we know the correct noise function  $\pi^*$  exactly, then maximizing  $L(\theta, \pi^*)$  over  $\theta$  is guaranteed to recover the correct

<sup>3</sup> We use the term “noise function” in order to be consistent with the terminology of [28] and to avoid deviating too much from the more standard terms “noise model” and “error model” used in the statistics and machine learning literature. We emphasize, however, that  $\pi$  defines much more than just the characteristics of experimental noise;  $\pi$  entirely specifies the relationship between measurements  $M$  and the underlying activity  $R$ . Were it not for prior terminology, the term “measurement function” might be preferable to “noise function.”



model  $\theta^*$ . However, if we assume an incorrect noise function  $\pi'$ , maximizing likelihood will typically recover an incorrect model  $\theta'$  (Fig. 4b).

### 4 Inference Using Mutual Information

Information theory provides an alternative inference approach. Suppose we hypothesize a specific model  $\theta$ , which gives predictions  $R$ . Denote the true model  $\theta^*$  and the corresponding true activity  $R^*$ . The dependence between  $S$ ,  $M$ ,  $R^*$ , and  $R$  will then form a Markov chain,

$$R \xleftarrow{\theta} S \xrightarrow{\theta^*} R^* \xrightarrow{\pi} M. \tag{5}$$

From the simple fact that  $M$  depends on  $R$  only through the value of  $R^*$ , any dependence measure  $\mathcal{D}$  that satisfies the data processing inequality (DPI) [29] must satisfy

$$\mathcal{D}[R; M] \leq \mathcal{D}[R^*; M]. \tag{6}$$

Therefore, in the set of possible models  $\theta$ , the true model is guaranteed to globally maximize the objective function  $\mathcal{D}(\theta) \equiv \mathcal{D}[R; M]$ .

One particularly relevant dependence measure that satisfies DPI is mutual information, a quantity that plays a fundamental role in information theory [29].<sup>4</sup> For the massively parallel experiments such as those in Fig. 2,  $R$  is continuous and  $M$  is discrete. In these cases, mutual information is given by

$$I(\theta) = I[R; M] = \sum_M \int dR p(R, M) \log \frac{p(R, M)}{p(R)p(M)}, \tag{7}$$

where  $p(M, R)$  is the joint distribution of activity predictions and measurements resulting from the model  $\theta$ . If one is able to estimate  $p(M, R)$  from a finite sample of data, mutual information can be used as an objective function for determining  $\theta$  without assuming any noise function  $\pi$ .

It should be noted that there are multiple dependence measures  $\mathcal{D}$  that satisfy DPI. One might wonder whether maximizing multiple different dependence measures would improve on the optimization of mutual information alone. The answer is not so simple. In [28] it was shown that if the correct model  $\theta^*$  is within the space of models under consideration, then, in the large data limit, maximizing mutual information is equivalent to simultaneously maximizing every dependence measure that satisfies DPI. On the other hand, one rarely has any assurance that the correct model  $\theta^*$  is within the space of parameterized models one is considering. In this case, considering different DPI-satisfying measures might provide a test for whether  $\theta^*$  is noticeably outside the space of parameterized models. To our knowledge, this potential approach to the model selection problem has yet to be demonstrated.

### 5 Relationship Between Likelihood and Mutual Information

A third inference approach is to admit that we do not know the noise function  $\pi$  a priori, and to fit both  $\theta$  and  $\pi$  simultaneously by maximizing  $L(\theta, \pi)$  over this pair. It is easy to see why this makes sense: the division of the inference problem into first measuring  $\pi$ , then learning  $\theta$  using that inferred  $\pi$ , is somewhat artificial. The process that maps  $S$  to  $M$  is determined

<sup>4</sup> See [43] for an extended discussion of mutual information as a measure of statistical association.

by both  $\theta$  and  $\pi$  and thus, from a probabilistic point of view, it makes sense to maximize likelihood over both of these quantities simultaneously.

We now show that, in the large  $N$  limit, maximizing likelihood over both  $\theta$  and  $\pi$  is equivalent to maximizing the mutual information between model predictions and measurements. Here we follow the argument given in [28]. In the large  $N$  limit, likelihood can be written

$$L(\theta, \pi) = \sum_M \int dR p(R, M) \log \pi(M|R) \quad (8)$$

$$= I(\theta) - D(\theta, \pi) - H[M], \quad (9)$$

where

$$D(\theta, \pi) = \sum_M \int dR p(R, M) \log \frac{p(M|R)}{\pi(M|R)}, \quad (10)$$

is the Kullback–Leibler divergence between the assumed noise function  $\pi$  and the observed noise function  $p(M|R)$ , and  $H[M] = -\sum_M p(M) \log p(M)$  is the entropy of the measurements, which does not depend on  $\theta$ . To maximize  $L(\theta, \pi)$  it therefore suffices to maximize  $I(\theta)$  over  $\theta$  alone, then to set the noise function  $\pi(M|R)$  equal to the empirical noise function  $p(M|R)$ , which causes  $D(\theta, \pi)$  to vanish.

Thus, when we are uncertain about the noise function  $\pi$ , we need not despair. We can, if we like, simply learn  $\pi$  at the same time that we learn  $\theta$ . We need not explicitly model  $\pi$  in order to do this; it suffices instead to maximize the mutual information  $I(\theta)$  over  $\theta$  alone.

The connection between mutual information and likelihood can further be seen in a quantity called the “noise-averaged” likelihood. This quantity was first described for the analysis of microarray data [27]; see also [28]. The central idea is to put an explicit prior on the space of possible noise functions, then compute likelihood after marginalizing over these noise functions. Explicitly, the per-datum log noise-averaged likelihood  $L_{\text{na}}(\theta)$  is related to  $L(\theta, \pi)$  via

$$e^{NL_{\text{na}}(\theta)} = \int d\pi p(\pi) e^{NL(\theta, \pi)}. \quad (11)$$

We will refer to  $L_{\text{na}}$  simply as “noise-averaged likelihood” in what follows.

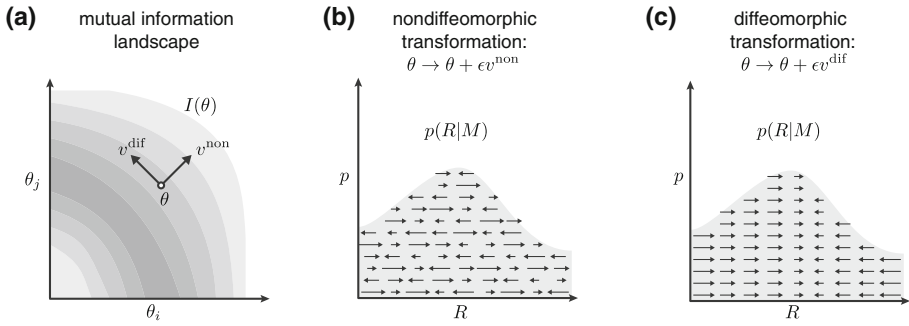
Under fairly general conditions, one finds that noise-averaged likelihood is related to mutual information via

$$L_{\text{na}}(\theta) = I(\theta) - \Delta(\theta) - H[M]. \quad (12)$$

Here, the effect of the noise function prior  $p(\pi)$  is absorbed entirely by the term  $\Delta(\theta)$ . Under weak assumptions,  $\Delta(\theta)$  vanishes in the  $N \rightarrow \infty$  limit and thus  $p(\pi)$  becomes irrelevant for the inference problem [27, 28].

## 6 Diffeomorphic Modes

Mutual information has a mathematical property that is important to account for when using it as an objective function: the mutual information between any two variables is unchanged by an invertible transformation of either variable. So if a change in model parameters,  $\theta \rightarrow \theta'$ , results in changes in model predictions  $R \rightarrow R'$  that preserves the rank order of these predictions, then



**Fig. 5** Illustration of diffeomorphic and nondiffeomorphic modes. **a** A diffeomorphic mode  $v^{\text{dif}}$  at a point  $\theta$  in parameter space is a vector that will (regardless of the underlying data) be tangent to a level curve of  $I(\theta)$ . All other vectors (e.g.,  $v^{\text{non}}$ ) correspond to nondiffeomorphic modes. **b** Moving  $\theta$  along a nondiffeomorphic mode results in a sort of “diffusion” in which the  $R$  values assigned to different sequences change rank order. Here, the probability distribution  $p(R|M)$  is illustrated (for fixed  $M$ ) in gray. The motion of individual  $R$  values upon such a change in  $\theta$  are indicated by arrows. **c** Changing  $\theta$  along a diffeomorphic mode, however, results in a “flow” of  $R$  values that maintains their rank order

$$I(\theta) = I[M; R] = I[M; R'] = I(\theta'), \tag{13}$$

and  $\theta$  and  $\theta'$  are judged to be equally valid.

By using mutual information as an objective function, we are therefore unable to constrain any parameters of  $\theta$  that, if changed, produce invertible transformations of model predictions. Such parameters are called “diffeomorphic parameters” or “diffeomorphic modes” [28]. The distinction between diffeomorphic modes and nondiffeomorphic modes is illustrated in Fig. 5.

### 6.1 Criterion for Diffeomorphic Modes

Following [28], we now derive a criterion that can be used to identify all of the diffeomorphic modes of a model  $\theta$ .<sup>5</sup> Consider an infinitesimal change in model parameters  $\theta \rightarrow \theta + d\theta$ , where the components of  $d\theta$  are specified by

$$d\theta_i = \epsilon v_i \tag{14}$$

for small  $\epsilon$  and for some vector  $v_i$  in  $\theta$ -space. This change in  $\theta$  will produce a corresponding change in model predictions  $R \rightarrow R + dR$ , where

$$dR = \epsilon \sum_i v_i \frac{\partial R}{\partial \theta_i}. \tag{15}$$

In general, the derivative  $\partial R / \partial \theta_i$  can have arbitrary dependence on the underlying sequence  $S$ . This transformation will preserve the rank order of  $R$ -values only if  $dR$  is the same for all sequences having the same value of  $R$ . The change  $dR$  must therefore be a function of  $R$  and have no other dependence on  $S$ . A diffeomorphic mode is a vector field  $v^{\text{dif}}(\theta)$  that has this property at all points in parameter space. Specifically, a vector field  $v^{\text{dif}}(\theta)$  is a diffeomorphic mode if and only if there is a function  $h(R, \theta)$  such that

$$\sum_i v_i^{\text{dif}}(\theta) \frac{\partial R}{\partial \theta_i} = h(R, \theta). \tag{16}$$

<sup>5</sup> Here, as throughout this paper, we restrict our attention to situations in which  $R$  is a scalar. The case of vector-valued model predictions  $R$  is worked out in [28].

### 6.2 Diffeomorphic Modes of Linear Models

As a simple example, consider a situation in which each sequence  $S$  is a  $D$ -dimensional vector and  $R$  is an affine function of  $S$ , i.e.

$$R = \theta_0 + \sum_{i=1}^D \theta_i S_i, \tag{17}$$

for model parameters  $\theta = \{\theta_0, \theta_1, \dots, \theta_D\}$ . The criterion in Eq. (16) then gives

$$v_0^{\text{dif}}(\theta) + \sum_{i=1}^D v_i^{\text{dif}}(\theta) S_i = h(R, \theta). \tag{18}$$

Because the left hand side is linear in  $S$  and  $R$  is linear in  $S$ , the function  $h(R, \theta)$  must be linear in  $R$ . Thus,  $h$  must have the form

$$h(R, \theta) = a(\theta) + b(\theta)R \tag{19}$$

for some functions  $a(\theta)$  and  $b(\theta)$ . The corresponding diffeomorphic mode is

$$v_i^{\text{dif}}(\theta) = \begin{cases} a(\theta) & i = 0 \\ b(\theta)\theta_i & i = 1, 2, \dots, D \end{cases}, \tag{20}$$

which has two degrees of freedom. Specifically, the  $a$  component of  $v^{\text{dif}}$  corresponds to adding a constant to  $R$  while the  $b$  component corresponds to multiplying  $R$  by a constant.

Note that if we had instead chosen  $R = \sum_{i=1}^D \theta_i S_i$ , i.e. left out the constant component  $\theta_0$ , then there would be only one diffeomorphic mode, corresponding to multiplication of  $R$  by a constant. This fact will be used when we analyze the Gaussian selection model in Sect. 8.

### 6.3 Diffeomorphic Modes of a Biophysical Model of Transcriptional Regulation

Diffeomorphic modes can become less trivial in more complicated situations. Consider the biophysical model of transcriptional regulation by the *E. coli lac* promoter (Fig. 3). This model was fit to Sort-Seq data in [12]. The form of this model is as follows. Let  $S$  denote a  $4 \times D$  matrix representing a DNA sequence of length  $D$  and having elements

$$S_{bl} = \begin{cases} 1 & \text{if base } b \text{ occurs at position } l \\ 0 & \text{otherwise} \end{cases} \tag{21}$$

where  $b \in \{A, C, G, T\}$  and  $l = 1, 2, \dots, D$ . The binding energy  $Q$  of CRP to DNA was modeled in [12] as an “energy matrix”: each position in the DNA sequence was assumed to contribute additively to the overall energy. Specifically,

$$Q = \sum_{b,l} \theta_Q^{bl} S_{bl} + \theta_Q^0, \tag{22}$$

where  $\theta_Q = \{\theta_Q^0, \theta_Q^{bl}\}$  are the parameters of this energy matrix. Similarly, the binding energy  $P$  of RNAP to DNA was modeled as

$$P = \sum_{b,l} \theta_P^{bl} S_{bl} + \theta_P^0. \tag{23}$$

Both energies were taken to be in thermal units ( $k_B T$ ). The rate of transcription  $R$  resulting from these binding energies was assumed to be proportional to the occupancy of RNAP at its binding site. This transcription rate is given by

$$R = R_{\max} \frac{e^{-P} + e^{-P-Q-\gamma}}{1 + e^{-Q} + e^{-P} + e^{-P-Q-\gamma}}, \tag{24}$$

where  $\gamma$  is the interaction energy between CRP and RNAP (again in units of  $k_B T$ ) and  $R_{\max}$  is a scalar.

Because the binding sites for CRP and RNAP do not overlap, one can learn the parameters  $\theta_Q$  and  $\theta_P$  from data separately by independently maximizing  $I[Q; M]$  and  $I[P; M]$ . Doing this, however, leaves undetermined the overall scale of each energy matrix as well as the chemical potentials  $\theta_P^0$  and  $\theta_Q^0$ . The reason is that the energy scale and chemical potential are diffeomorphic modes of energy matrix models and therefore cannot be inferred by maximizing mutual information.

However, if  $Q$  and  $P$  are inferred together by maximizing  $I[R; M]$  instead, one is now able to learn both energy matrices with a physically meaningful energy scale. The chemical potential of CRP,  $\theta_Q^0$ , is also determined. The only parameters left unspecified are the chemical potential of RNA polymerase,  $\theta_P^0$ , and the maximal transcription rate  $R_{\max}$ . The reason for this is that in the formula for  $R$  in Eq. (24) the energies  $P$  and  $Q$  combine in a nonlinear way. This nonlinearity eliminates three of the four diffeomorphic modes of  $P$  and  $Q$ .<sup>6</sup> See [28] for the derivation of this result.

### 6.4 Dual Modes of the Noise Function

Diffeomorphic transformations of model parameters can be thought of as being equivalent to certain transformations of the noise function. Consider the transformation of model parameters

$$\theta_i \rightarrow \theta'_i = \theta_i + \epsilon v_i, \tag{25}$$

where  $\epsilon$  is an infinitesimal number and  $v_i$  is a vector in  $\theta$ -space.<sup>7</sup> For any sequence  $S$ , this transformation induces a transformation of the model prediction

$$R \rightarrow R' = R + \epsilon \sum_i v_i \frac{\partial R}{\partial \theta_i}. \tag{26}$$

To see the effect this transformation has on likelihood, we rewrite Eq. (4) as,

$$L(\theta, \pi) = \langle \log \pi(M|R) \rangle_{\text{data}}, \tag{27}$$

where  $\langle \cdot \rangle_{\text{data}}$  indicates an average taken over the measurements  $M_n$  and predictions  $R_n$  for all of the sequences  $S_n$  in the data set. The change in likelihood resulting from Eq. (26) is therefore given by

$$L(\theta', \pi) = L(\theta, \pi) + \epsilon \left\langle \frac{\partial \log \pi(M|R)}{\partial R} \sum_i \frac{\partial R}{\partial \theta_i} v_i \right\rangle_{\text{data}}. \tag{28}$$

<sup>6</sup> The one additional diffeomorphic mode is created by the introduction of the parameter  $R_{\max}$ .

<sup>7</sup> For the sake of clarity we suppress the  $\theta$ -dependence of  $v^{\text{dif}}$ ,  $\tilde{v}^{\text{dif}}$ , and  $h(R)$  in what follows.

Now suppose that there is a noise function  $\pi'$  that has an equivalent effect on likelihood, i.e.,

$$L(\theta', \pi) = L(\theta, \pi') + O(\epsilon^2), \tag{29}$$

for all possible data sets  $\{S_n, M_n\}$ . We say that this transformation of the noise function  $\pi \rightarrow \pi'$  is “dual” to the transformation  $\theta \rightarrow \theta'$  of model parameters. The transformed noise function will necessarily have the form

$$\log \pi'(M|R) = \log \pi(M|R) + \epsilon \tilde{v}(M, R) \tag{30}$$

for some function  $\tilde{v}(M, R)$ . To determine  $\tilde{v}$  we consider the transformation of likelihood induced by  $\pi \rightarrow \pi'$ :

$$L(\theta, \pi') = L(\theta, \pi) + \epsilon \langle \tilde{v}(M, R) \rangle_{\text{data}}. \tag{31}$$

Comparing Eqs. (28) and (31), we see that  $\pi \rightarrow \pi'$  will be dual to  $\theta \rightarrow \theta'$  for all possible data sets if and only if

$$\frac{\partial \log \pi(M|R)}{\partial R} \sum_i \frac{\partial R}{\partial \theta_i} v_i = \tilde{v}(M, R) \tag{32}$$

for all sequences  $S$ .

For general choice of vector  $v$ , no function  $\tilde{v}$  will exist that satisfies Eq. (32). The reason is that  $\partial R/\partial \theta_i$  will typically depend on the sequence  $S$  independently of the value of  $R$ . In other words, for a fixed value of  $M$  and  $R$ , the left hand side of Eq. (32) will retain a dependence on  $S$ . The right hand side, however, cannot have such a dependence. The converse is also true: for general choice of the function  $\tilde{v}$ , no vector  $v$  will exist such that Eq. (32) is satisfied for all sequences. This is evident from the simple fact that  $v$  is a finite dimensional vector while  $\tilde{v}$  is a function of the continuous quantity  $R$  and therefore has an infinite number of degrees of freedom.

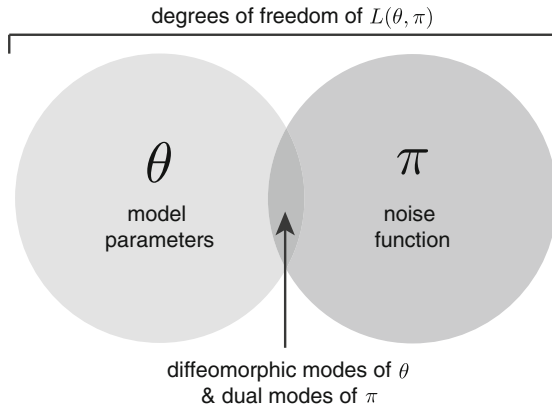
In fact, Eq. (32) will have a solution if and only if

$$\sum_i \frac{\partial R}{\partial \theta_i} v_i^{\text{dif}} = h(R) \tag{33}$$

for some function  $h$ . Here we have added the superscript “dif” because this is precisely the definition of a diffeomorphic mode given in Eq. (16). In this case, the function  $\tilde{v}^{\text{dif}}$  dual to this diffeomorphic mode  $v^{\text{dif}}$  is seen to be

$$\tilde{v}^{\text{dif}}(M, R) = \frac{\partial \log \pi(M|R)}{\partial R} h(R). \tag{34}$$

These findings are summarized by the Venn diagram in Fig. 6. Arbitrary transformations of the model parameters  $\theta$  will alter likelihood in a way that cannot be imitated by any change to the noise function  $\pi$ . The reverse is also true: most changes to  $\pi$  cannot be imitated by a corresponding change in  $\theta$ . However, a subset of transformations of  $\theta$  are equivalent to corresponding dual transformations of  $\pi$ . These transformations are precisely the diffeomorphic transformations of  $\theta$ . This partial duality between  $\theta$  and  $\pi$  has a simple interpretation: the choice of how we parse an experiment into an activity model  $\theta$  and a noise function  $\pi$  is not unique. The ambiguity in this choice is parameterized by the diffeomorphic modes of  $\theta$  and the dual modes of  $\pi$ .



**Fig. 6** Venn diagram illustrating the degrees of freedom of the likelihood  $L(\theta, \pi)$  considered over all possible data sets  $\{S_n, M_n\}$ . Altering the model parameters  $\theta$  will typically change  $L(\theta, \pi)$  in a way that cannot be recapitulated by changes in the noise function  $\pi$ . Similarly, changes in  $\pi$  cannot typically be imitated by changes in  $\theta$ . However, diffeomorphic transformations of  $\theta$  will affect  $L(\theta, \pi)$  in the exact same way that dual transformation of  $\pi$  will. The diffeomorphic modes of  $\theta$  and the dual modes of  $\pi$  can therefore be thought of as lying within the intersection of  $\theta$  and  $\pi$

### 7 Error Bars from Likelihood, Mutual Information, and Noise-Averaged Likelihood

We now consider the consequences of performing inference using various objective functions at large but finite  $N$ . Specifically, we discuss the optimal parameters and corresponding error bars that are found by sampling  $\theta$  from posterior distributions of the form

$$p(\theta|\text{data}) \sim e^{NF(\theta)} \tag{35}$$

for the following choices of the objective function  $F(\theta)$ :

- (a)  $F(\theta) = L(\theta, \pi^*)$  is likelihood computed using the correct noise function  $\pi^*$ .
- (b)  $F(\theta) = L(\theta, \pi')$  where  $\pi'$  differs from  $\pi^*$  by a small but arbitrary error.
- (c)  $F(\theta) = L(\theta, \pi'')$  where  $\pi''$  differs from  $\pi^*$  by a small amount along a dual mode.
- (d)  $F(\theta) = I(\theta)$  is the mutual information between measurements and model predictions.
- (e)  $F(\theta) = L_{\text{na}}(\theta)$  is the noise-averaged likelihood.

To streamline notation, we will use  $\langle \cdot \rangle$  to denote averages computed in multiple different contexts. In each case, the appropriate context will be specified by a subscript. As above  $\langle \cdot \rangle_{\text{data}}$  will denote averaging over a specific data set  $\{S_n, M_n\}_{n=1}^N$ .  $\langle \cdot \rangle_{\text{real}}$  will indicate averaging over an infinite number of data set realizations.  $\langle \cdot \rangle_S$ ,  $\langle \cdot \rangle_{S,M}$ ,  $\langle \cdot \rangle_{S|R}$ , and  $\langle \cdot \rangle_{S|R,M}$  will respectively denote averages over the distributions  $p(S)$ ,  $p(S, M)$ ,  $p(S|R)$ , and  $p(S|R, M)$ , the empirical distributions obtained in the infinite data limit.  $\langle \cdot \rangle_{\theta}$  will indicate an average computed over parameter values  $\theta$  sampled from the posterior distribution  $p(\theta|\text{data})$ . Subscripts on  $\text{cov}(\cdot)$  or  $\text{var}(\cdot)$  should be interpreted analogously.

#### 7.1 Likelihood

Consider Eq.(35) with  $F(\theta) = L(\theta, \pi^*)$  at large but finite  $N$ . The posterior distribution  $p(\theta|\text{data})$  will, in general, be maximized at some choice of parameters  $\theta^o$  that deviates

randomly from the correct parameters  $\theta^*$ . At large  $N$ ,  $p(\theta|\text{data})$  will become sharply peaked about  $\theta^o$  with a peak width governed by the Hessian of likelihood; specifically

$$\text{cov}_\theta(\theta_i - \theta_i^o, \theta_j - \theta_j^o) = -\frac{H_{ij}^{-1}}{N}, \tag{36}$$

where

$$H_{ij} = \left. \frac{\partial^2 L(\theta, \pi^*)}{\partial \theta_i \partial \theta_j} \right|_{\theta^*}, \tag{37}$$

is the Hessian of the likelihood. It is also readily shown (see Appendix 1) that this peak width is consistent with the correct parameters  $\theta^*$ , in the sense that

$$\text{cov}_{\text{real}}(\theta_i^* - \theta_i^o, \theta_j^* - \theta_j^o) = \text{cov}_\theta(\theta_i - \theta_i^o, \theta_j - \theta_j^o). \tag{38}$$

In Appendix 1 we show that the Hessian of likelihood, Eq. (37), is given by

$$H_{ij} = - \int dR p(R) J(R) \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R} \bigg|_{\theta^*}, \tag{39}$$

where

$$J(R) = \sum_M \pi^*(M|R) \left[ \frac{\partial \log \pi^*(M|R)}{\partial R} \right]^2 = - \sum_M \pi^*(M|R) \frac{\partial^2 \log \pi^*(M|R)}{\partial R^2} \tag{40}$$

is the Fisher information of the noise function  $\pi^*$ . This Fisher information is a nonnegative measure of how sensitive our experiment is in the vicinity of  $R$ .<sup>8</sup> We thus see that, as long as the set of vectors  $\partial R/\partial \theta_i$  spans all directions in parameter space, the Hessian matrix  $H_{ij}$  will be nonsingular. Using  $F(\theta) = L(\theta, \pi^*)$  will therefore put constraints on all directions in parameter space, and these constraints will shrink with increasing data as  $N^{-1/2}$ . This situation is illustrated in Fig. 7a.

Now consider what happens if instead we use a noise function  $\pi'$  that deviates from  $\pi^*$  in a small but arbitrary way. Specifically, let

$$\log \pi'(M|R) = \log \pi^*(M|R) + \epsilon f(M, R) \tag{41}$$

for some function  $f(M, R)$  and small parameter  $\epsilon$ . It is readily shown (see Appendix 1) that the maximum likelihood parameters  $\theta'$  will deviate from  $\theta^*$  by an amount

$$\langle \theta'_i - \theta_i^* \rangle_{\text{real}} = -\epsilon \sum_j H_{ij}^{-1} w_j, \quad \text{where} \quad w_j = \left\langle \frac{\partial f}{\partial R} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|\theta^*}. \tag{42}$$

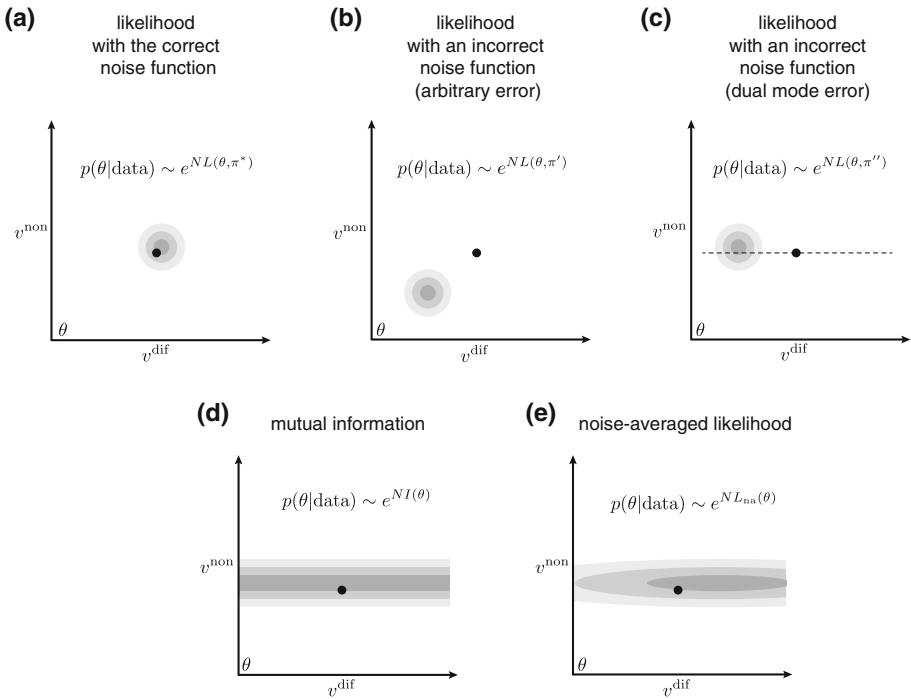
This expected deviation does not depend on  $N$  and will therefore not shrink to zero in the large  $N$  limit. Indeed, for any choice of  $\epsilon > 0$ , there will always be an  $N$  large enough such that this bias in  $\theta'$  dominates over the uncertainty due to finite sampling.

Is there any restriction on the types of biases in  $\theta'$  that can be produced by the choice of incorrect noise function  $\pi'$ ? In general, no. Because the Hessian matrix  $H$  is nonsingular, one can always find a vector  $w$  such that the deviation of  $\theta'$  from  $\theta^*$  in Eq. (42) points in any chosen direction of  $\theta$ -space. As long as the functions

$$g_i(R) = \left\langle \frac{\partial R}{\partial \theta_i} \right\rangle_{S|R} \bigg|_{\theta^*} \tag{43}$$

<sup>8</sup> In what follows we assume that  $J(R) > 0$  almost everywhere. This just reflects the assumption that our experiment actually does convey information about  $R$  through the measurements  $M$  it provides.





**Fig. 7** Posterior distributions on model parameters resulting from various objective functions. Each panel schematically illustrates the posterior distribution  $p(\theta|\text{data})$  (gray shaded area) as it relates to the correct model  $\theta^*$  (dot) along both diffeomorphic (abscissa) and nondiffeomorphic (ordinate) directions in parameter space. **a** Likelihood with the correct noise function  $\pi^*$  leads to a posterior distribution consistent with  $\theta^*$  in all parameters. **b** Likelihood with a noise function  $\pi'$  that differs arbitrarily from  $\pi^*$  will, in general, lead to a posterior distribution that is inconsistent with  $\theta^*$  along both diffeomorphic and nondiffeomorphic modes. **c** Likelihood with a noise function  $\pi''$  that differs from  $\pi^*$  only along a dual mode  $\tilde{v}^{\text{dif}}$  leads to a posterior that is inconsistent with  $\theta^*$  only along the diffeomorphic mode  $v^{\text{dif}}$  (parallel to dashed line), but consistent with  $\theta^*$  in all other directions (perpendicular to dashed line). **d** Using mutual information gives a posterior that is consistent with  $\theta^*$ ; this posterior places constraints similar to likelihood along non-diffeomorphic modes but places no constraints whatsoever along diffeomorphic modes. **e** Using noise-averaged likelihood results in a posterior distribution similar to mutual information but with weak constraints on diffeomorphic modes resulting from the noise function prior  $p(\pi)$

are linearly independent for different indices  $i$ , a function  $f$  can always be found that generates the vector  $w$  in Eq. (42).

We therefore see that arbitrary errors in the noise function will bias the inference of model parameters in arbitrary directions. This fact presents a major concern for standard likelihood-based inference: if you assume an incorrect noise function  $\pi$ , the parameters  $\theta$  that you then infer will, in general, be biased in an unpredictable way. Moreover, the magnitude of this bias will be directly proportional to the magnitude of the error in the log of your assumed noise function. This problem is illustrated in Fig. 7b.

There is a case that deserves some additional consideration. Suppose we use a noise function  $\pi''$  that differs from  $\pi^*$  only along a dual mode  $\tilde{v}^{\text{dif}}$ , i.e.,

$$\log \pi''(M|R) = \log \pi^*(M|R) + \epsilon \tilde{v}^{\text{dif}}(M, R). \tag{44}$$

The maximum likelihood parameters  $\theta''$  of  $L(\theta, \pi'')$  will still deviate from  $\theta^*$  by an amount that does not shrink to zero in the  $N \rightarrow \infty$  limit. However, this bias in parameter values will be restricted to the diffeomorphic mode  $v^{\text{dif}}$  to which  $\tilde{v}^{\text{dif}}$  is dual, i.e.,

$$\langle \theta_i'' - \theta_i^* \rangle_{\text{real}} = -\epsilon v_i^{\text{dif}}. \tag{45}$$

This state of affairs ain't so bad since the incorrect noise function will lead to model parameters that are inaccurate only along modes that we already know we cannot learn from the data. This situation is illustrated in Fig. 7c; see Appendix 1 for the derivation of Eq. (45).

### 7.2 Mutual Information

The constraints on parameters imposed by using mutual information  $I(\theta)$  as the objective function  $F(\theta)$  in Eq. (35) are determined by the Hessian

$$K_{ij} = \left. \frac{\partial^2 I(\theta)}{\partial \theta_i \partial \theta_j} \right|_{\theta^*}. \tag{46}$$

Appendix 2 provides a detailed derivation of this Hessian, which after some computation is found to be given by

$$K_{ij} = - \int dR p(R) J(R) \left[ \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R} - \left\langle \frac{\partial R}{\partial \theta_i} \right\rangle_{S|R} \left\langle \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R} \right] \Big|_{\theta^*}. \tag{47}$$

Comparing Eqs. (47) and (39), we see that for any vector  $v$  in parameter space,

$$- \sum_{i,j} H_{ij} v_i v_j \geq - \sum_{i,j} K_{ij} v_i v_j \geq 0. \tag{48}$$

Likelihood is thus seen to constrain parameters in all directions at least as much as mutual information does. As expected, mutual information provides no constraint whatsoever in the direction of any diffeomorphic mode  $v^{\text{dif}}$  of the model, since

$$- \sum_{i,j} K_{ij} v_i^{\text{dif}} v_j^{\text{dif}} = \int dR p(R) J(R) \left[ \langle h^2(R) \rangle_{S|R} - \langle h(R) \rangle_{S|R}^2 \right] \Big|_{\theta^*} = 0. \tag{49}$$

The converse is also true: if there is no constraint on parameters along  $v$ , then  $v$  must be a diffeomorphic mode. This is because

$$- \sum_{i,j} K_{ij} v_i v_j = \int dR p(R) J(R) \text{ var} \left( \sum_i v_i \frac{\partial R}{\partial \theta_i} \right) \Big|_{S|R, \theta^*}. \tag{50}$$

Because  $J(R)$  is positive almost everywhere, the right hand side of Eq. (50) can vanish only if  $\sum_i v_i \frac{\partial R}{\partial \theta_i}$  does not differ between any two sequences that have the same  $R$  value. There must therefore exist a function  $h(R)$  such that  $h(R) = \sum_i v_i \frac{\partial R}{\partial \theta_i}$  for all sequences  $S$ . This is precisely the requirement in Eq. (16) that  $v$  be a diffeomorphic mode.

However, except along diffeomorphic modes, we can generally expect that the constraints provided by likelihood and by mutual information will be of the same magnitude. This situation is illustrated in Fig. 7d. Indeed, in the next section we will see an explicit example where all nondiffeomorphic constraints imposed by mutual information are comparable to those imposed by likelihood.

Before proceeding, we note that the relationship between the Hessians of likelihood and mutual information suggests an analogy to fluid mechanics. Consider a trajectory in parameter space given by  $\theta_i(t) = tv_i$ , where  $t$  is time and  $v$  is a velocity vector pointing in the direction of motion. This motion in parameter space will induce a motion in the prediction  $R(t)$  that the model provides for every sequence  $S$ . The set of sequences  $\{S_n\}$  thus presents us with a dynamic cloud of “particles” moving about in  $R$ -space. At  $t = 0$ , the quantity  $\langle \dot{R}^2 \rangle_{S|R}$  will be proportional to the average kinetic energy of particles at location  $R$ . The quantity  $\langle \dot{R} \rangle_{S|R}^2$  will be proportional to the (per particle) kinetic energy of the bulk fluid element at  $R$ , a quantity that does not count energy due to thermal motion. In this way we see that  $-\sum_{i,j} H_{ij} v_i v_j$  is a weighted tally of total kinetic energy, whereas  $-\sum_{i,j} K_{ij} v_i v_j$  corresponds to a tally of internal thermal energy only, the kinetic energy of bulk motion having been subtracted out.

### 7.3 Noise-Averaged Likelihood

Noise-averaged likelihood provides constraints in between those of likelihood, computed using the correct noise function, and those of mutual information. This is illustrated in Fig. 7e. Whereas mutual information provides no constraints whatsoever on the diffeomorphic modes of  $\theta$ , noise-averaged likelihood provides weak constraints in these directions. These soft constraints reflect the Hessian of  $\Delta(\theta)$  in Eq. (12). The constraints along diffeomorphic modes, however, have an upper bound on how tight they can become in the  $N \rightarrow \infty$  limit. This is because such constraints only reflect our prior  $p(\pi)$  on the noise function, not the information we glean from data.

## 8 Worked Example: Gaussian Selection

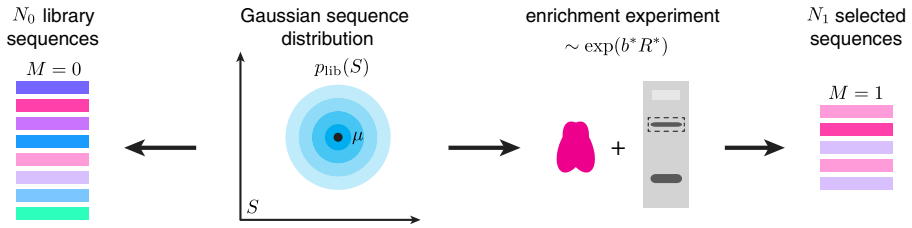
The above principles can be illustrated in the following analytically tractable model of a massively parallel experiment, which we call the “Gaussian selection model.” In this model, our experiment starts with a large library of “DNA” sequences  $S$ , each of which is actually a  $D$ -dimensional vector drawn from a Gaussian probability distribution<sup>9</sup>

$$p_{\text{lib}}(S) = (2\pi)^{-D/2} \exp\left(-\frac{|S - \mu|^2}{2}\right). \tag{51}$$

Here,  $\mu$  is a  $D$ -dimensional vector defining the average sequence in the library. From this library we extract sequences into two bins, labeled  $M = 0$  and  $M = 1$ . We fill the  $M = 0$  bin with sequences sampled indiscriminately from the library. The  $M = 1$  bin is filled with sequences sampled from this library with relative probability

$$\frac{p(M = 1|S)}{p(M = 0|S)} = \exp(a^* + b^* R^*) \tag{52}$$

<sup>9</sup> For the sake of simplicity we set the covariance matrix of this distribution equal to the identity matrix. The more general case of a non-identity covariance matrix yields the same basic results. Also, we note that, while approximating discrete DNA sequences by continuous vectors might seem crude, it is only the marginal distributions  $p(R|M)$  that matter for the inference problem. Most of the quantities  $R$  that one encounters in practice are computed by summing up contributions from a large number of different nucleotide positions. In such cases, the marginal distributions  $p(R|M)$  will often be nearly continuous and virtually indistinguishable from the marginal distributions one might obtain from a Gaussian sequence library.



**Fig. 8** Illustration of the Gaussian selection model of a massively parallel experiment. Each assayed sequence in this model is a  $D$ -dimensional vector. The library (corresponding to bin  $M = 0$ ) consists of  $N_0$  sequences  $S$  drawn from a Gaussian distribution  $p_{\text{lib}}(S)$  that is centered on a specific sequence  $\mu$ . Bin  $M = 1$  consists of  $N_1$  sequences drawn from the distribution  $p_{\text{lib}}(S)$  then enriched by a factor of  $\exp(b^* R^*)$  where  $R^* = S^T \theta^*$ . This enrichment procedure is analogous to selecting protein-bound DNA sequences where  $b^* R^*$  is negative the binding energy. Calculations in the text are performed in the  $N_0 \gg N_1$  limit

where the activity  $R^*$  is defined as the dot product of  $S$  with a  $D$ -dimensional vector  $\theta^*$ , i.e.,

$$R^* = S^T \theta^*. \tag{53}$$

We use  $N_M$  to denote the number of sequences in each bin  $M$ , along with  $N = N_0 + N_1$ .

All of our calculations are performed in the limit where  $N_1$  is large but  $N_0$  is far larger. More specifically, we assume that  $\exp(a^* + b^* R^*) \ll 1$  everywhere that both  $p(S|M = 0)$  and  $p(S|M = 1)$  are significant. We use  $\epsilon$  to denote the ratio

$$\epsilon \equiv \frac{p(M = 1)}{p(M = 0)} = \frac{N_1}{N_0}, \tag{54}$$

and all of our calculations are carried out only to first order in  $\epsilon$ . This model experiment is illustrated in Fig. 8.

Our goal is this: given the sampled sequences in the two bins, recover the parameters  $\theta^*$  defining the sequence–function relationship for  $R^*$ . To do this, we adopt the following model for the sequence-dependent activity  $R$ :

$$R = S^T \theta, \tag{55}$$

where  $\theta$  is the  $D$ -dimensional vector we wish to infer. From the arguments above and in [28], it is readily seen that the magnitude of  $\theta$ , i.e.  $|\theta|$ , is the only diffeomorphic mode of the model: changing this parameter rescales  $R$ , preserving rank order.

### 8.1 Bin-Specific Distributions

We can readily calculate the conditional sequence distribution  $p(S|M)$  for each bin  $M$ , as well as the conditional distribution  $p(R|M)$  of model predictions. Because the sequences sampled for bin 0 are indiscriminately drawn from  $p_{\text{lib}}$ , we have

$$p(S|M = 0) = p_{\text{lib}}(S) = (2\pi)^{-D/2} \exp\left(-\frac{|S - \mu|^2}{2}\right). \tag{56}$$

The distribution of selected sequences is found to be

$$p(S|M = 1) = (2\pi)^{-D/2} \exp\left(-\frac{|S - \mu - b^* \theta^*|^2}{2}\right). \tag{57}$$

The value of  $\epsilon$  is found to be related to  $a^*$ ,  $b^*$ , and  $\theta^*$  via

$$\epsilon = \exp\left(a^* + b^* \mu^T \theta^* + \frac{b^{*2} |\theta^*|^2}{2}\right). \quad (58)$$

Appendix 3 provides an explicit derivation of Eqs. (57) and (58).

We compute the distribution of model predictions for each bin as follows. For each bin  $M$ , this distribution is defined as

$$p(R|M) = \int dS \delta(R - \theta^T S) p(S|M). \quad (59)$$

This can be analytically calculated for both of the bins owing to the Gaussian form of each sequence distribution. We find that

$$p(R|M=0) = \frac{1}{\sqrt{2\pi}|\theta|} \exp\left(-\frac{(R - \mu^T \theta)^2}{2|\theta|^2}\right), \quad (60)$$

$$p(R|M=1) = \frac{1}{\sqrt{2\pi}|\theta|} \exp\left(-\frac{(R - [\mu + b^* \theta^*]^T \theta)^2}{2|\theta|^2}\right). \quad (61)$$

See Appendix 3 for details.

## 8.2 Noise Function

To compute likelihood, we must posit a noise function  $\pi(M|R)$ . Based on our prior knowledge of the selection procedure, we choose  $\pi(M|R)$  so that

$$\frac{\pi(M=1|R)}{\pi(M=0|R)} = \exp(a + bR), \quad (62)$$

where  $a$  and  $b$  are scalar parameters that we might or might not know *a priori*. This, combined with the normalization requirement,  $\sum_M \pi(M|R) = 1$ , gives

$$\pi(M=1|R) = \frac{e^{a+bR}}{1 + e^{a+bR}}, \quad \pi(M=0|R) = \frac{1}{1 + e^{a+bR}}. \quad (63)$$

This noise function  $\pi$  is correct when  $a = a^*$  and  $b = b^*$ . The parameter  $b$  is dual to the diffeomorphic mode  $|\theta|$ , whereas the parameter  $a$  is not dual to any diffeomorphic mode.

In the experimental setup used to motivate the Gaussian selection model, the parameter  $a$  is affected by many aspects of the experiment, including the concentration of the protein used in the binding assay, the efficiency of DNA extraction from the gel, and the relative amount of PCR amplification used for the bin 0 and bin 1 sequences. In practice, these aspects of the experiment are very hard to control, much less predict. From the results in the previous section, we can expect that if we assume a specific value for  $a$  and perform likelihood-based inference, inaccuracies in this value for  $a$  will distort our inferred model  $\theta$  in an unpredictable (i.e., nondiffeomorphic) way. We will, in fact, see that this is the case. The solution to this problem, of course, is to infer  $\theta$  alone by maximizing the mutual information  $I(\theta)$ ; in this case the values for  $a$  and  $b$  become irrelevant. Alternatively, one can place a prior on  $a$  and  $b$ , then maximize noise-averaged likelihood  $L_{\text{na}}(\theta)$ . We now analytically explore the consequences of these three approaches.

### 8.3 Likelihood

Using the noise function in Eq.(63), the likelihood  $L$  becomes a function of  $\theta$ ,  $a$ , and  $b$ . Computing  $L$  in the  $N \rightarrow \infty$  and  $\epsilon \rightarrow 0$  limits, we find that

$$L(\theta, a, b) = \epsilon[a + b\theta^T \mu + bb^* \theta^T \theta^*] - \exp\left(a + b\theta^T \mu + \frac{b^2|\theta|^2}{2}\right). \tag{64}$$

We now consider the consequences of various approaches for using  $L(\theta, a, b)$  to estimate  $\theta^*$ . In each case, the inferred optimum will be denoted by a superscript ‘o’. Standard likelihood-based inference requires that we assume a specific value for  $a$  and for  $b$ , then optimize  $L(\theta, a, b)$  over  $\theta$  alone by setting

$$0 = \frac{\partial L}{\partial \theta_i} \Big|_{\theta^o, a, b} \tag{65}$$

for each component  $i$ . By this criteria we find that the optimal model  $\theta^o$  is given by a linear combination of  $\theta^*$  and  $\mu$ :

$$\theta^o = \frac{cb^*}{b}\theta^* + \frac{c-1}{b}\mu, \tag{66}$$

where  $c$  is a scalar that solves the transcendental equation

$$c = \exp\left([a^* - a] + \frac{1-c^2}{2}|b^*\theta^* + \mu|^2\right). \tag{67}$$

See Appendix 2 for the derivation of this result. Note that  $c$  is determined only by the value of  $a$  and not by the value of  $b$ . Moreover,  $c = 1$  if and only if  $a = a^*$ .

If our assumed noise function is correct, i.e.,  $a = a^*$  and  $b = b^*$ , then

$$\theta^o = \theta^*. \tag{68}$$

Thus, maximizing likelihood will identify the correct model parameters. This exemplifies the general behavior illustrated in Fig. 7a.

If  $a = a^*$  but  $b \neq b^*$ , our assumed noise function will differ from the correct noise function only in a manner dual to the diffeomorphic mode  $|\theta|$ . In this case we find that  $c = 1$  and

$$\theta^o = \frac{b^*}{b}\theta^*. \tag{69}$$

$\theta^o$  is thus proportional but not equal to  $\theta^*$ . This comports with our claim above that the diffeomorphic mode of the inferred model, i.e.  $|\theta^o|$ , will be biased so as to compensate for the error in the dual parameter  $b$ . This finding follows the behavior described in Fig. 7c.

If  $a \neq a^*$ , however,  $c \neq 1$ . As a result,  $\theta^o$  is a nontrivial linear combination of  $\theta^*$  and  $\mu$ , and will thus point in a different direction than  $\theta^*$ . This is true regardless of the value of  $b$ . This behavior is illustrated in Fig. 7b: errors in non-dual parameters of the noise function will typically lead to errors in nondiffeomorphic parameters of the activity model.

We now consider the error bars that likelihood places on model parameters. Setting  $\theta = \theta^o + \delta\theta$  and expanding  $L(\theta, a, b)$  about  $\theta^o$ , we find that

$$NL(\theta, a^*, b^*) \approx NL(\theta^o, a^*, b^*) - \frac{N_1 b^{*2}}{2} \sum_{i,j} \Lambda_{ij} \delta\theta_i \delta\theta_j, \tag{70}$$

where  $\Lambda_{ij} = \delta_{ij} + (\mu_i + b^*\theta_i^*)(\mu_j + b^*\theta_j^*)$ . Note that all eigenvalues of  $\Lambda$  are greater or equal to 1. Adopting the posterior distribution

$$p(\theta|\text{data}) \sim e^{NL(\theta,a,b)} \tag{71}$$

therefore gives a covariance matrix on  $\theta$  of

$$\langle \delta\theta_i \delta\theta_j \rangle = \frac{\Lambda_{ij}^{-1}}{N_1 b^{*2}}. \tag{72}$$

Thus,  $\delta\theta \sim N_1^{-1/2}$  in all directions of  $\theta$ -space. Therefore, when the noise function is incorrect and  $N$  is sufficiently large, the finite bias introduced into  $\theta^o$  will cause  $\theta^*$  to fall outside the inferred error bars.

### 8.4 Mutual Information

In the  $\epsilon \rightarrow 0$  limit, Eq. (7) simplifies to

$$I(\theta) = \epsilon \int dR p(R|M = 1) \log \frac{p(R|M = 1)}{p(R|M = 0)} + O(\epsilon^2). \tag{73}$$

The lowest order term on the right hand side can be evaluated exactly using Eqs. (60) and (61):

$$I(\theta) = \frac{\epsilon b^{*2}}{2} \frac{(\theta^T \theta^*)^2}{|\theta|^2}. \tag{74}$$

See Appendix 3 for details. Note that the expression on the right is invariant under a rescaling of  $\theta$ . This reflects the fact that  $|\theta|$  is a diffeomorphic mode of the model defined in Eq. (55).

To find the model  $\theta^o$  that maximizes mutual information, we set

$$0 = \left. \frac{\partial I}{\partial \theta_i} \right|_{\theta^o} = \frac{\epsilon b^{*2} \theta^{oT} \theta^*}{|\theta^o|^2} \left[ \theta_i^* - \theta_i^o \frac{\theta^{oT} \theta^*}{|\theta^o|^2} \right] \tag{75}$$

The optimal model  $\theta^o$  must therefore be parallel to  $\theta^*$ , i.e.

$$\theta^o \propto \theta^*. \tag{76}$$

Expanding about  $\theta = \theta^o + \delta\theta$  as above, we find that

$$NI(\theta) = NI(\theta^o) - \frac{N_1 b^{*2}}{2} (\delta\theta_{\perp})^2 \tag{77}$$

where  $\delta\theta_{\perp}$  is the component of  $\delta\theta$  perpendicular to  $\theta^*$ ; see Appendix 3. Therefore, if we use the posterior distribution  $p(\theta|\text{data}) \sim e^{NI(\theta)}$  to infer  $\theta$ , we find uncertainties in directions perpendicular to  $\theta^*$  of magnitude  $N_1^{-1/2}$ . These error bars are only slightly larger than those obtained using likelihood, and have the same dependence on  $N$ . However, we find no constraint whatsoever on the component of  $\delta\theta$  parallel to  $\theta^*$ . These results are illustrated by Fig. 7d.

## 8.5 Noise-Averaged Likelihood

We can also compute the noise-averaged likelihood,  $L_{\text{na}}(\theta)$ , in the case of a uniform prior on  $a$  and  $b$ , i.e.  $p(\pi) = p(a, b) = \mathcal{C}$  where  $\mathcal{C}$  is an infinitesimal constant. We find that

$$\exp[NL_{\text{na}}(\theta)] = \int d\pi p(\pi) \exp[NL(\theta, \pi)] \quad (78)$$

$$= \mathcal{C} \int_{-\infty}^{\infty} da \int_{-\infty}^{\infty} db \exp\left(N_1[a + b\theta^T \mu + bb^* \theta^T \theta^*] - N \exp\left[a + b\theta^T \mu + \frac{b^2 |\theta|^2}{2}\right]\right) \quad (79)$$

$$= \mathcal{C} \Gamma(N_1) \sqrt{\frac{2\pi}{N|\theta|^2}} \exp\left(\frac{N_1 b^{*2} (\theta^T \theta^*)^2}{2 |\theta|^2}\right). \quad (80)$$

See the Appendix 3 for details. Thus,

$$L_{\text{na}}(\theta) = I(\theta) - \frac{1}{N} \log |\theta| + \text{const}, \quad (81)$$

where the constant (which absorbs  $\mathcal{C}$  entirely) does not depend on  $\theta$ . If we perform Bayesian inference using noise-averaged likelihood, i.e., using  $p(\theta|\text{data}) \sim e^{NL_{\text{na}}(\theta)}$ , we will therefore find in the large  $N$  limit that  $\delta\theta_{\perp}$  is constrained in the same way as if we had used mutual information. The noise function prior we have assumed further results in weak constraints on  $|\theta|$  that do not tighten as  $N$  increases.<sup>10</sup> This is represented in Fig. 7e.

## 9 Discussion

The systematic study of quantitative sequence–function relationships in biology is just now becoming possible thanks to the development of a variety of massively parallel experiments. Concepts and methods from statistical physics are likely to prove valuable for understanding this basic class of biological phenomena as well as for learning sequence–function relationships from data.

In this paper we have discussed the problem of learning parametric models of sequence–function relationships from experiments having poorly characterized experimental noise. We have seen that standard likelihood-based inference, which requires an explicit model of experimental noise, will generally lead to incorrect model parameters due to errors in the assumed noise function. By contrast, mutual-information-based inference allows one to learn parametric models without having to assume any noise function at all. Mutual-information-based inference is unable to pin down the values of model parameters along diffeomorphic modes. This behavior reflects a fundamental difference between how diffeomorphic and non-diffeomorphic modes are constrained by data. Diffeomorphic modes arise from arbitrariness in the distinction between the activity model and the noise function. These findings were illustrated using an analytically tractable model for a massively parallel experiment.

The study of quantitative sequence–function relationships still presents many challenges, both theoretical and computational. One major practical difficulty with the mutual-information-based approach described here is accurately estimating mutual information

<sup>10</sup> In the case at hand,  $|\theta^{\circ}|$  is pushed all the way to zero. This is an artifact of the simple flat prior  $p(a, b)$ . If we instead adopt a weak Gaussian prior on  $b$ , we can still carry out the computation of  $L_{\text{na}}$  analytically, and in this case we find that  $|\theta^{\circ}|$  is finite.



from data. Although methods are available for doing this [44], it remains unclear whether any are accurate enough to enable computational sampling of the posterior distribution  $p(\theta|\text{data}) \sim e^{NI(\theta)}$ , as suggested here. Moreover, none of these estimation methods is regarded as definitive. We believe this lack of clarity regarding how to estimate mutual information reflects the fact that the density estimation problem itself has never been fully solved, even in one or two dimensions. We are hopeful, however, that field-theoretic methods for estimating probability densities [45–47] might help resolve the problem of mutual information.

The problem of model selection poses a major theoretical challenge. Ideally, one would like to explore a hierarchy of possible model classes when fitting parametric models to data. However, when considering effective models it is unclear how to move far beyond independent site models (e.g., energy matrices) due to the number of parameters growing exponentially with the length of the sequence. Moreover, when learning mechanistic models such as the model of the *lac* promoter featured in Fig. 3, it is unclear how to go about systematically testing different arrangements of binding sites, different protein–protein interactions, and so on. We emphasize that this model prioritization problem is fundamentally theoretical, not computational, and as of now there is little clarity on how to address this matter.

Finally, the geometric structure of sequence–function relationships presents an array of intriguing questions. For instance, very little is known (in any system) about how convex or glassy such landscapes in sequence space are, what their density of states looks like, etc.. Most of the biological and evolutionary implications of these aspects of sequence–function relationships also have yet to be worked out. We believe that the methods and ideas of statistical physics may lead to important insights into these questions in the near future.

**Acknowledgments** We would like to thank L. Peliti, O. Revoire, and T. Mora for organizing this special issue. This work was supported by the Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory and the Starr Cancer Consortium (17-A723).

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## Appendix 1: Maximum Likelihood Under Various Noise Functions

At the correct noise function  $\pi^*$ , likelihood is given by

$$L(\theta, \pi^*) = \langle \log \pi^*(M|R) \rangle_{\text{data}}. \quad (82)$$

Taylor expanding this quantity about  $\theta^*$  gives

$$\begin{aligned} L(\theta, \pi^*) &= L(\theta^*, \pi^*) + \sum_i \frac{\partial L}{\partial \theta_i} \Big|_{\theta^*} (\theta_i - \theta_i^*) \\ &+ \frac{1}{2} \sum_{i,j} \frac{\partial^2 L}{\partial \theta_i \partial \theta_j} \Big|_{\theta^*} (\theta_i - \theta_i^*)(\theta_j - \theta_j^*) + \dots \end{aligned} \quad (83)$$

We define the random vector  $u$  in terms of the coefficient of the linear term of this expansion:

$$\frac{u_i}{\sqrt{N}} \equiv \frac{\partial L}{\partial \theta_i} \Big|_{\theta^*} = \left\langle \frac{\partial \log \pi^*(M|R)}{\partial R} \frac{\partial R}{\partial \theta_i} \right\rangle_{\text{data}} \Big|_{\theta^*}. \quad (84)$$

Because  $u_i/\sqrt{N}$  is defined as a sum of  $N$  random terms, and because the mean of these terms vanishes, the covariance  $\langle u_i u_j \rangle_{\text{real}}$  will, by the central limit theorem, be given by

$$\langle u_i u_j \rangle_{\text{real}} = \left\langle \left[ \frac{\partial \log \pi^*(M|R)}{\partial R} \right]^2 \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S,M} \Big|_{\theta^*} \tag{85}$$

$$= \sum_M \int dR p(M, R) \left[ \frac{\partial \log \pi^*(M|R)}{\partial R} \right]^2 \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R, M} \Big|_{\theta^*}. \tag{86}$$

At  $\theta = \theta^*$ , each measurement  $M$  will provide no additional information about  $S$  beyond that provided by the model prediction  $R = \theta(S)$ . Mathematically this means that

$$p(S|R, M)|_{\theta^*} = p(S|R)|_{\theta^*} \tag{87}$$

for all  $S, R$ , and  $M$ . Equivalently, the conditional expectation value of any sequence-dependent function  $f(S)$  will obey

$$\langle f(S) \rangle_{S|R, M} \Big|_{\theta^*} = \langle f(S) \rangle_{S|R} \Big|_{\theta^*} \tag{88}$$

for all  $M$ . We use this fact to simplify Eq. (86):

$$\langle u_i u_j \rangle_{\text{real}} = \int dR p(R) \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R} \sum_M \pi(M|R) \left[ \frac{\partial \log \pi^*(M|R)}{\partial R} \right]^2 \Big|_{\theta^*} \tag{89}$$

$$= \int dR p(R) J(R) \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R} \Big|_{\theta^*} \tag{90}$$

where  $J(R)$  is the Fisher information from Eq. (40).

We compute the Hessian of likelihood as follows:

$$H_{ij} \equiv \frac{\partial^2 L(\theta, \pi^*)}{\partial \theta_i \partial \theta_j} \Big|_{\theta^*} = \left\langle \frac{\partial^2 \log \pi^*(M|R)}{\partial R^2} \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S,M} \Big|_{\theta^*} + \left\langle \frac{\partial \log \pi^*(M|R)}{\partial R} \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S,M} \Big|_{\theta^*}. \tag{91}$$

The second term on the right hand side vanishes because of Eq. (88):

$$\left\langle \frac{\partial \log \pi^*(M|R)}{\partial R} \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S,M} \Big|_{\theta^*} = \sum_M \int dR p(R, M) \frac{\partial \log \pi^*(M|R)}{\partial R} \left\langle \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S|R, M} \Big|_{\theta^*} \tag{92}$$

$$= \int dR p(R) \left\langle \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S|R} \left[ \sum_M \pi(M|R) \frac{\partial \log \pi^*(M|R)}{\partial R} \right] \Big|_{\theta^*} \tag{93}$$

$$= \int dR p(R) \left\langle \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S|R} \left[ \frac{\partial}{\partial R} \sum_M \pi(M|R) \right] \Big|_{\theta^*} \tag{94}$$

$$= \int dR p(R) \left\langle \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S|R} \left[ \frac{\partial}{\partial R} 1 \right] \Big|_{\theta^*} \tag{95}$$

$$= 0. \tag{96}$$

We therefore find that

$$H_{ij} = \sum_M \int dR p(R, M) \frac{\partial^2 \log \pi^*(M|R)}{R^2} \Big|_{\theta^*} \tag{97}$$

$$= - \int dR p(R) J(R) \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} \right\rangle_{S|R} \Big|_{\theta^*}, \tag{98}$$

which is Eq. (39). Note that, from Eq. (90),  $\langle u_i u_j \rangle_{\text{real}} = -H_{ij}$ .

The optimum  $\theta^o$  of  $L(\theta, \pi^*)$  will occur when

$$0 = \frac{\partial L(\theta, \pi^*)}{\partial \theta} \Big|_{\theta^o} = \frac{u_i}{\sqrt{N}} + \sum_j H_{ij} (\theta_i^o - \theta_j^*) + \dots \tag{99}$$

We therefore find that, to lowest order in  $N^{-1/2}$ ,

$$\theta_i^o = \theta_i^* - \sum_j H_{ij}^{-1} \frac{u_j}{\sqrt{N}}. \tag{100}$$

The covariance of  $\theta^o$  is thus given by

$$\left\langle (\theta_i^o - \theta_i^*) (\theta_j^o - \theta_j^*) \right\rangle_{\text{real}} = \sum_{k,l} H_{ik}^{-1} \frac{\langle u_k u_l \rangle_{\text{real}}}{N} H_{lj}^{-1} = -\frac{H_{ij}^{-1}}{N}, \tag{101}$$

which is Eq. (38).

Under the incorrect noise function  $\pi'$  defined in Eq. (41),

$$L(\theta, \pi') = L(\theta, \pi^*) + \epsilon \langle f(M, R) \rangle_{\text{data}} \tag{102}$$

$$\approx L(\theta^*, \pi^*) + \epsilon \langle f(M, R) \rangle_{S|\theta^*} + \sum_i \left[ \frac{u_i}{\sqrt{N}} + \epsilon w_i \right] (\theta_i - \theta_i^*) + \frac{1}{2} \sum_{ij} H_{ij} (\theta_i - \theta_i^*) (\theta_j - \theta_j^*) + \dots \tag{103}$$

where

$$w_i = \left\langle \frac{\partial f}{\partial R} \frac{\partial R}{\partial \theta_i} \right\rangle_{S|\theta^*}. \tag{104}$$

Let  $\theta'$  denote the maximum of  $L(\theta, \pi')$ . Setting  $\frac{\partial L(\theta, \pi')}{\partial \theta_i} = 0 \Big|_{\theta'}$ , we find

$$\theta'_i = \theta_i^* - \sum_j H_{ij}^{-1} \left[ \frac{u_j}{\sqrt{N}} + \epsilon w_j \right], \tag{105}$$

from which we get Eq. (42).

In the case of a noise function  $\pi''$  that differs from  $\pi^*$  only along a dual mode, as in Eq. (44), the vector  $w_i$  is given by

$$w_i = \left\langle \frac{\partial \tilde{v}^{\text{dif}}}{\partial R} \frac{\partial R}{\partial \theta_i} \right\rangle_{S|\theta^*}. \tag{106}$$

The maximum likelihood parameters  $\theta''$  will therefore satisfy

$$0 = \sum_j H_{ij} \left\langle \theta_j'' - \theta_j^* \right\rangle_{\text{real}} + \epsilon w_i \tag{107}$$

$$= \sum_j \left\langle \frac{\partial^2 \log \pi}{\partial R^2} \frac{\partial R}{\partial \theta_i} \frac{\partial R}{\partial \theta_j} + \frac{\partial \log \pi}{\partial R} \frac{\partial^2 R}{\partial \theta_i \partial \theta_j} \right\rangle_{S, M} \Big|_{\theta=\theta^*} \left\langle \theta_j'' - \theta_j^* \right\rangle_{\text{real}} + \epsilon \left\langle \frac{\partial R}{\partial \theta_i} \frac{\partial}{\partial R} \left[ \frac{\partial \log \pi}{\partial R} h(R) \right] \right\rangle_{S, M} \Big|_{\theta=\theta^*} \tag{108}$$

$$= \left\langle \frac{\partial}{\partial \theta_i} \frac{\partial \log \pi}{\partial R} \left[ \sum_j \frac{\partial R}{\partial \theta_j} \left\langle \theta_j'' - \theta_j^* \right\rangle_{\text{real}} + \epsilon h(R) \right] \right\rangle_{S, M} \Big|_{\theta=\theta^*} \tag{109}$$

$$= \left\langle \frac{\partial}{\partial \theta_i} \frac{\partial \log \pi}{\partial R} \sum_j \frac{\partial R}{\partial \theta_j} \left( \left\langle \theta_j'' - \theta_j^* \right\rangle_{\text{real}} + \epsilon v_j^{\text{dif}} \right) \right\rangle_{S, M} \Big|_{\theta=\theta^*}, \tag{110}$$

which is solved by Eq. (45). The fact that this uniquely specifies  $\langle \theta_i'' - \theta_i^* \rangle_{\text{real}}$  follows from the Hessian  $H$  being nonsingular.

### Appendix 2: Gradient and Hessian of Mutual Information

Here we calculate the gradient and Hessian of mutual information evaluated at  $\theta = \theta^*$ . We do this by first computing derivatives of the empirical probability distributions  $p(R)$  and  $p(R, M)$  with respect to model parameters. The mathematical trick used to do this is adapted from [31]. These results are first applied to likelihood in order to demonstrate their use and correctness. We then use this approach to compute the gradient and Hessian of mutual information. To clarify these derivations, we use  $r(\theta, S)$ , instead of  $\theta(S)$ , to explicitly denote the model prediction  $R$  as a function of sequence  $S$  and model parameters  $\theta$ . We also define  $\partial_i \equiv \frac{\partial}{\partial \theta_i}$  and use  $\int dS$  to represent sums over sequences.

#### How the Distribution of Model Predictions Changes with Model Parameters

The empirical probability distribution of model predictions  $R$  is given by

$$p(R) = \int dS p(S) \delta(R - r(\theta, S)). \tag{111}$$

The gradient of this probability distribution with respect to model parameters is computed as follows:

$$\partial_i p(R) = \int dS p(S) \partial_i \delta(R - r(\theta, S)) \tag{112}$$

$$= - \int dS p(S) \left[ \frac{\partial}{\partial R} \delta(R - r(\theta, S)) \right] \partial_i r \tag{113}$$

$$= - \frac{\partial}{\partial R} \left[ p(R) \int dS p(S|R) \delta(R - r(\theta, S)) \partial_i r \right] \tag{114}$$

$$= - \frac{\partial}{\partial R} \left[ p(R) \langle \partial_i r \rangle_{S|R} \right]. \tag{115}$$

Similarly, the Hessian of  $p(R)$  is given by

$$\partial_i \partial_j p(R) = \int dS p(S) \left\{ \left[ \frac{\partial^2}{\partial R^2} \delta(R - r(\theta, S)) \right] \partial_i r \partial_j r - \left[ \frac{\partial}{\partial R} \delta(R - r(\theta, S)) \right] \partial_i \partial_j r \right\} \tag{116}$$

$$= \frac{\partial^2}{\partial R^2} \left[ p(R) \langle \partial_i r \partial_j r \rangle_{S|R} \right] - \frac{\partial}{\partial R} \left[ p(R) \langle \partial_i \partial_j r \rangle_{S|R} \right]. \tag{117}$$

Analogous results follow for the gradient and Hessian of the joint distribution  $p(R, M)$ :

$$\partial_i p(R, M) = - \frac{\partial}{\partial R} \left[ p(R, M) \langle \partial_i r \rangle_{S|R, M} \right], \tag{118}$$

$$\partial_i \partial_j p(R, M) = \frac{\partial^2}{\partial R^2} \left[ p(R, M) \langle \partial_i r \partial_j r \rangle_{S|R, M} \right] - \frac{\partial}{\partial R} \left[ p(R, M) \langle \partial_i \partial_j r \rangle_{S|R, M} \right]. \tag{119}$$

### Gradient and Hessian of Likelihood

Likelihood can be expressed in terms of the empirical distribution  $p(R, M)$  as

$$L(\theta, \pi) = \sum_M \int dR p(R, M) \log \pi(M|R). \tag{120}$$

Keep in mind that  $R$  is just a dummy variable in this integral; the empirical distribution  $p$  is the only quantity that depends on  $\theta$ . The gradient of likelihood is therefore computed as

$$\partial_i L = \sum_M \int dR [\partial_i p(R, M)] \log \pi(M|R) \tag{121}$$

$$= \sum_M \int dR \left\{ - \frac{\partial}{\partial R} \left[ p(R, M) \langle \partial_i r \rangle_{S|R, M} \right] \right\} \log \pi(M|R) \tag{122}$$

$$= \sum_M \int dR p(R, M) \frac{\partial \log \pi(M|R)}{\partial R} \langle \partial_i r \rangle_{S|R, M} \tag{123}$$

$$= \left\langle \frac{\partial \log \pi(M|R)}{\partial R} \partial_i r \right\rangle_{S, M}. \tag{124}$$

Note that in going from Eqs. (122) to (123) we used integration by parts. The Hessian of likelihood is computed similarly:

$$\partial_i \partial_j L = \sum_M \int dR [\partial_i \partial_j p(R, M)] \log \pi(M|R) \tag{125}$$

$$= \sum_M \int dR \log \pi(M|R) \left\{ \frac{\partial^2}{\partial R^2} \left[ p(R, M) \langle \partial_i r \partial_j r \rangle_{S|R, M} \right] - \frac{\partial}{\partial R} \left[ p(R, M) \langle \partial_i \partial_j r \rangle_{S|R, M} \right] \right\} \tag{126}$$

$$\begin{aligned}
 &= \sum_M \int dR p(R, M) \left\{ \frac{\partial^2 \log \pi(M|R)}{\partial R^2} \langle \partial_i r \partial_j r \rangle_{S|R, M} \right. \\
 &\quad \left. + \frac{\partial \log \pi(M|R)}{\partial R} \langle \partial_i \partial_j r \rangle_{S|R, M} \right\}. \tag{127}
 \end{aligned}$$

This expression is valid for all choices  $\theta$  and  $\pi$ .

Restricting our attention now to  $\theta = \theta^*$  and  $\pi = \pi^*$ , we see that the second term in Eq. (127) vanishes as it did in Eq. (92) through Eq. (96). Moreover, the first term gives

$$\partial_i \partial_j L = - \int dR p(R) J(R) \langle \partial_i r \partial_j r \rangle_{S|R}, \tag{128}$$

which is the formula obtained for  $H_{ij}$  in Eq. (39).

### Gradient and Hessian of Mutual Information

The gradient and Hessian computations for mutual information are simplified by expressing mutual information in terms of its component entropies. We write

$$I(\theta) = H_R(\theta) + H_M - H_{RM}(\theta) \tag{129}$$

where

$$H_{RM}(\theta) = - \sum_M \int dR p(R, M) \log p(M, R), \tag{130}$$

$$H_R(\theta) = - \int dR p(R) \log p(R), \tag{131}$$

$$H_M = - \sum_M p(M) \log p(M). \tag{132}$$

The gradient of  $H_R$  is given by

$$\partial_i H_R = - \int dR [\partial_i p(R)] \log p(R) - \int dR p(R) \partial_i \log p(R) \tag{133}$$

$$= - \int dR [\partial_i p(R)] \log p(R) - \int dR p(R) \frac{1}{p(R)} \partial_i p(R) \tag{134}$$

$$= - \int dR [\partial_i p(R)] \log p(R) - \partial_i 1 \tag{135}$$

$$= - \int dR [\partial_i p(R)] \log p(R). \tag{136}$$

Similarly,

$$\partial_i H_{RM} = - \sum_M \int dR [\partial_i p(R, M)] \log p(R, M). \tag{137}$$

$H_M$  does not depend on  $\theta$ , so  $\partial_i H_M = 0$ . The resulting gradient of mutual information is

$$\partial_i I = \sum_M \int dR [\partial_i p(R, M)] \log p(R, M) - \int dR [\partial_i p(R)] \log p(R) \tag{138}$$

$$= \sum_M \int dR [\partial_i p(R, M)] \log \frac{p(R, M)}{p(R)} \tag{139}$$

$$= \sum_M \int dR [\partial_i p(R, M)] \log p(M|R). \tag{140}$$

Note from Eq. (121) that  $\partial_i I = \partial_i L$  whenever  $\pi(M|R) = p(M|R)$ .

Now let's compute the Hessian of  $H_R$ :

$$\partial_i \partial_j H_R = - \int dR [\partial_i \partial_j p(R)] \log p(R) - \int dR [\partial_i p(R)] \partial_j \log p(R) \tag{141}$$

$$= - \int dR [\partial_i \partial_j p(R)] \log p(R) - \int dR p(R) [\partial_i \log p(R)] [\partial_j \log p(R)]. \tag{142}$$

Similarly,

$$\begin{aligned} \partial_i \partial_j H_{RM} &= - \sum_M \int dR [\partial_i \partial_j p(R, M)] \log p(R, M) \\ &\quad - \sum_M \int dR p(R, M) [\partial_i \log p(R, M)] [\partial_j \log p(R, M)]. \end{aligned} \tag{143}$$

The Hessian of mutual information is therefore given by,

$$\partial_i \partial_j I = \partial_i \partial_j H_R - \partial_i \partial_j H_{RM}. \tag{144}$$

Using the form of  $\partial_i \partial_j L$  in Eq. (125), we see that this reduces to

$$\partial_i \partial_j I = \partial_i \partial_j L + \Lambda_{ij}^{RM} - \Lambda_{ij}^R, \tag{145}$$

where

$$\Lambda_{ij}^R = \int dR p(R) [\partial_i \log p(R)] [\partial_j \log p(R)] \tag{146}$$

and

$$\Lambda_{ij}^{RM} = \sum_M \int dR p(R, M) [\partial_i \log p(R, M)] [\partial_j \log p(R, M)]. \tag{147}$$

We now split  $\Lambda_{ij}^R$  and  $\Lambda_{ij}^{RM}$  into four terms each. For  $\Lambda_{ij}^R$  we get

$$\Lambda_{ij}^R = \int dR p(R) \left\{ -\frac{1}{p(R)} \frac{\partial}{\partial R} [p(R) \langle \partial_i r \rangle_{S|R}] \right\} \left\{ -\frac{1}{p(R)} \frac{\partial}{\partial R} [p(R) \langle \partial_j r \rangle_{S|R}] \right\} \tag{148}$$

$$\begin{aligned} &= \int dR p(R) \left\{ \frac{\partial \log p(R)}{\partial R} \langle \partial_i r \rangle_{S|R} + \frac{\partial}{\partial R} \langle \partial_i r \rangle_{S|R} \right\} \\ &\quad \times \left\{ \frac{\partial \log p(R)}{\partial R} \langle \partial_j r \rangle_{S|R} + \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \right\} \end{aligned} \tag{149}$$

$$= A_{ij}^R + B_{ij}^R + B_{ji}^R + C_{ij}^R, \tag{150}$$

where

$$A_{ij}^R = \int dR p(R) \left[ \frac{\partial \log p(R)}{\partial R} \right]^2 \langle \partial_i r \rangle_{S|R} \langle \partial_j r \rangle_{S|R}, \tag{151}$$

$$B_{ij}^R = \int dR p(R) \left[ \frac{\partial \log p(R)}{\partial R} \right] \langle \partial_i r \rangle_{S|R} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R}, \tag{152}$$

$$C_{ij}^R = \int dR p(R) \left[ \frac{\partial}{\partial R} \langle \partial_i r \rangle_{S|R} \right] \left[ \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \right]. \tag{153}$$

Similarly,

$$\Lambda_{ij}^{RM} = A_{ij}^{RM} + B_{ij}^{RM} + B_{ji}^{RM} + C_{ij}^{RM} \tag{154}$$

where

$$A_{ij}^{RM} = \sum_M \int dR p(R, M) \left[ \frac{\partial \log p(R, M)}{\partial R} \right]^2 \langle \partial_i r \rangle_{S|R, M} \langle \partial_j r \rangle_{S|R, M}, \tag{155}$$

$$B_{ij}^{RM} = \sum_M \int dR p(R, M) \left[ \frac{\partial \log p(R, M)}{\partial R} \right] \langle \partial_i r \rangle_{S|R, M} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R, M}, \tag{156}$$

$$C_{ij}^{RM} = \sum_M \int dR p(R, M) \left[ \frac{\partial}{\partial R} \langle \partial_i r \rangle_{S|R, M} \right] \left[ \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R, M} \right]. \tag{157}$$

It is unclear how to simplify the expression for  $\partial_i \partial_j I$  at general choices of  $\theta$ . At  $\theta = \theta^*$ , however, the expectation value  $\langle \partial_i r \rangle_{S|R, M}$  loses all  $M$ -dependence and this causes a lot of cancellations to occur:

$$C_{ij}^{RM} = \sum_M \int dR p(R, M) \left[ \frac{\partial}{\partial R} \langle \partial_i r \rangle_{S|R} \right] \left[ \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \right] \tag{158}$$

$$= \int dR p(R) \left[ \frac{\partial}{\partial R} \langle \partial_i r \rangle_{S|R} \right] \left[ \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \right] \tag{159}$$

$$= C_{ij}^R \tag{160}$$

and

$$B_{ij}^{RM} = \int dR p(R) \left[ \sum_M p(M|R) \frac{\partial \log p(R, M)}{\partial R} \right] \langle \partial_i r \rangle_{S|R} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \tag{161}$$

$$= \int dR p(R) \left[ \sum_M \frac{p(M|R)}{p(R, M)} \frac{\partial p(R, M)}{\partial R} \right] \langle \partial_i r \rangle_{S|R} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \tag{162}$$

$$= \int dR p(R) \left[ \frac{1}{p(R)} \frac{\partial}{\partial R} \sum_M p(R, M) \right] \langle \partial_i r \rangle_{S|R} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \tag{163}$$

$$= \int dR p(R) \left[ \frac{1}{p(R)} \frac{\partial p(R)}{\partial R} \right] \langle \partial_i r \rangle_{S|R} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \tag{164}$$

$$= \int dR p(R) \left[ \frac{\partial \log p(R)}{\partial R} \right] \langle \partial_i r \rangle_{S|R} \frac{\partial}{\partial R} \langle \partial_j r \rangle_{S|R} \tag{165}$$

$$= B_{ij}^R. \tag{166}$$



We therefore find that,

$$\Lambda_{ij}^{RM} - \Lambda_{ij}^R = A_{RM} - A_R \tag{167}$$

$$= \int dR p(R) \langle \partial_i r \rangle_{S|R} \langle \partial_j r \rangle_{S|R} \times \left\{ \sum_M p(M|R) \left[ \frac{\partial \log p(R, M)}{\partial R} \right]^2 - \left[ \frac{\partial \log p(R)}{\partial R} \right]^2 \right\}. \tag{168}$$

The expression in braces can be simplified as follows:

$$\begin{aligned} & \sum_M p(M|R) \left[ \frac{\partial \log p(R, M)}{\partial R} \right]^2 - \left[ \frac{\partial \log p(R)}{\partial R} \right]^2 \\ &= \sum_M p(M|R) \left\{ \left[ \frac{\partial \log p(M|R)}{\partial R} + \frac{\partial \log p(R)}{\partial R} \right]^2 - \left[ \frac{\partial \log p(R)}{\partial R} \right]^2 \right\} \end{aligned} \tag{169}$$

$$= \sum_M p(M|R) \left\{ \left[ \frac{\partial \log p(M|R)}{\partial R} \right]^2 + \frac{\partial \log p(R)}{\partial R} \frac{\partial \log p(M|R)}{\partial R} \right\} \tag{170}$$

$$= J(R) + \frac{1}{p(R)} \frac{\partial p(R)}{\partial R} \frac{\partial}{\partial R} \sum_M p(M|R) \tag{171}$$

$$= J(R). \tag{172}$$

The Hessian of mutual information at  $\theta = \theta^*$  therefore has a rather simple form:

$$K_{ij} = H_{ij} + \Lambda_{ij}^{RM} - \Lambda_{ij}^R = - \int dR p(R) J(R) \left[ \langle \partial_i r \partial_j r \rangle_{S|R} - \langle \partial_i r \rangle_{S|R} \langle \partial_j r \rangle_{S|R} \right], \tag{173}$$

which is Eq. (47).

### Appendix 3: Gaussian Selection Model

#### Derivation of Eqs. (57) and (58)

Applying Bayes’s theorem twice,

$$p(S|M = 1) = \frac{p(M = 1|S)}{p(M = 1)} p(S) = \frac{p(M = 1|S)}{p(M = 1)} \frac{p(M = 0)}{p(M = 0|S)} p(S|M = 0). \tag{174}$$

Using Eqs. (56), (52), and (54) then gives

$$p(S|M = 1) = \epsilon^{-1} e^{a^* + b^* S^T \theta^*} (2\pi)^{-D/2} \exp\left(-\frac{|S - \mu|^2}{2}\right). \tag{175}$$

Next we complete the square in the exponent:

$$-\frac{|S - \mu|^2}{2} + b^* S^T \theta^* = -\frac{|S|^2 + |\mu|^2 - 2\mu^T S - 2b^* S^T \theta^*}{2} \tag{176}$$

$$= -\frac{|S|^2 + |\mu|^2 + |b^* \theta^*|^2 - 2\mu^T S - 2b^* S^T \theta^* + 2b^* \mu^T \theta^*}{2} + \frac{|b^* \theta^*|^2}{2} + b^* \mu^T \theta^* \tag{177}$$

$$= -\frac{|S - \mu - b^* \theta^*|^2}{2} + \frac{|b^* \theta^*|^2}{2} + b^* \mu^T \theta^*. \tag{178}$$

From the first term in Eq. (178) we recover Eq. (57). To get  $\epsilon$ , we substitute Eq. (178) into Eq. (175). Comparing this to Eq. (57) then gives

$$1 = \epsilon^{-1} e^{a^*} \exp\left(\frac{|b^* \theta^*|^2}{2} + b^* \mu^T \theta^*\right). \tag{179}$$

Solving for  $\epsilon$  recovers Eq. (58).

**Derivation of Eqs. (60) and (61)**

Here we describe how to compute  $p(R|M)$  where  $R = \theta^T S$ . We first consider the case of  $M = 0$ .

$$p(R|M = 0) = \int dS p(S|M = 0) \delta(R - S^T \theta) \tag{180}$$

$$= \int dS p(S|M = 0) \delta([R - \mu^T \theta] - [S - \mu]^T \theta) \tag{181}$$

$$= \int dS p(S|M = 0) \delta(R' - S'^T \theta) \tag{182}$$

where  $R' = R - \mu^T \theta$  and  $S' = S - \mu$ . We have chosen to work with  $R'$  and  $S'$  instead of  $R$  and  $S$  because  $p(S'|M = 0)$  is centered about 0. Now, split  $S'$  up into the components parallel and perpendicular to  $\theta$ :

$$S' = S'_\perp + S'_\parallel \hat{\theta}, \tag{183}$$

where  $S'_\perp$  is a vector orthogonal to  $\theta$ ,  $S'_\parallel$  is a scalar, and  $\hat{\theta} = \theta/|\theta|$ . This definition gives  $S'^T \theta = S'_\parallel |\theta|$ . Continuing with the integration,

$$p(R|M = 0) = \int dS'_\perp \int_{-\infty}^{\infty} dS'_\parallel \delta(R' - S'_\parallel |\theta|) (2\pi)^{-D/2} \exp\left(-\frac{S'^2_\perp}{2} - \frac{S'^2_\parallel}{2}\right) \tag{184}$$

$$= \int_{-\infty}^{\infty} dS'_\parallel \delta(R' - S'_\parallel |\theta|) (2\pi)^{-1/2} \exp\left(-\frac{S'^2_\parallel}{2}\right) \tag{185}$$

$$= \int_{-\infty}^{\infty} dS'_\parallel \delta\left(\frac{R'}{|\theta|} - S'_\parallel\right) |\theta|^{-1} (2\pi)^{-1/2} \exp\left(-\frac{S'^2_\parallel}{2}\right) \tag{186}$$

$$= |\theta|^{-1} (2\pi)^{-1/2} \exp\left(-\frac{R'^2}{2|\theta|^2}\right). \tag{187}$$

Finally, substituting  $R$  back for  $R'$  gives

$$p(R|M = 0) = \frac{1}{\sqrt{2\pi}|\theta|} \exp\left(-\frac{(R - \mu^T\theta)^2}{2|\theta|^2}\right). \tag{188}$$

To compute  $p(R|M = 1)$ , we just replace  $\mu \rightarrow \mu + b^*\theta^*$ , giving

$$p(R|M = 1) = \frac{1}{\sqrt{2\pi}|\theta|} \exp\left(-\frac{(R - [\mu + b^*\theta^*]^T\theta)^2}{2|\theta|^2}\right). \tag{189}$$

**Derivation of Eq. (64)**

We compute likelihood in the  $N \rightarrow \infty$  limit as follows:

$$L(\theta, a, b) = \sum_M p(M) \int dR p(R|M) \log \pi(M|R) \tag{190}$$

$$= \frac{N_0}{N} \int dR p(R|M = 0) \log \frac{1}{1 + e^{a+bR}} + \frac{N_1}{N} \int dR p(R|M = 1) \log \frac{e^{a+bR}}{1 + e^{a+bR}} \tag{191}$$

$$\approx -\frac{N_0}{N} \int dR p(R|M = 0) e^{a+bR} + \frac{N_1}{N} \int dR p(R|M = 1) [a + bR] \tag{192}$$

$$\approx -\langle e^{a+bR} \rangle_{S|M=0} + \epsilon \langle a + bR \rangle_{S|M=1}. \tag{193}$$

In deriving Eq.(193) we assumed that  $e^{a+bR} \ll 1$  for all values of  $R$  over which both  $p(R|M = 0)$  and  $p(R|M = 1)$  have significant support. This assumption necessarily holds in the  $\epsilon \rightarrow 0$  limit. We have also kept only the lowest order terms in  $\epsilon$ . Note in particular that  $\langle e^{a+bR} \rangle_{S|M=0}$  will be of order  $\epsilon$ .

The second term in Eq.(193) can be directly read off from Eq. (61):

$$\langle a + bR \rangle_{S|M=1} = a + b \langle R \rangle_{S|M=1} = a + b\mu^T\theta + bb^*\theta^T\theta^*. \tag{194}$$

From Eq.(60) we see that the first term in Eq.(193) can be computed by completing the square:

$$-\frac{(R - \mu^T\theta)^2}{2|\theta|^2} + bR = -\frac{R^2 + (\mu^T\theta)^2 - 2(\mu^T\theta)R - 2b|\theta|^2R}{2|\theta|^2} \tag{195}$$

$$= -\frac{R^2 + (\mu^T\theta)^2 + b^2|\theta|^4 - 2(\mu^T\theta)R - 2b|\theta|^2R + 2(\mu^T\theta)b|\theta|^2}{2|\theta|^2} + b(\mu^T\theta) + \frac{b^2|\theta|^2}{2} \tag{196}$$

$$= -\frac{(R - \mu^T\theta - b|\theta|^2)^2}{2|\theta|^2} + b(\mu^T\theta) + \frac{b^2|\theta|^2}{2}, \tag{197}$$

from which we get

$$\left\langle e^{a+bR} \right\rangle_{S|M=1} = \exp \left[ a + b(\mu^T \theta) + \frac{b^2|\theta|^2}{2} \right]. \tag{198}$$

Plugging Eqs. (194) and (198) into Eq. (193) gives the formula for  $L(\theta, a, b)$  in Eq. (64).

**Derivation of Eqs. (66) and (67)**

Here we show how to derive the optimal  $\theta$  for  $L(\theta, a, b)$ , with  $a$  and  $b$  fixed. Setting the gradient of  $L$  with respect to  $\theta$  to zero,

$$0 = \frac{\partial L}{\partial \theta_i} \Big|_{\theta^o} = \epsilon b(\mu_i + b^* \theta_i^*) - b(\mu_i + b \theta_i^o) \exp \left( a + b \mu^T \theta^o + \frac{b^2 |\theta^o|^2}{2} \right). \tag{199}$$

This gives

$$\mu_i + b \theta_i^o = \epsilon (\mu_i + b^* \theta_i^*) \exp \left( -a - b \mu^T \theta^o - \frac{b^2 |\theta^o|^2}{2} \right) \tag{200}$$

$$= c (\mu_i + b^* \theta_i^*) \tag{201}$$

where  $c$  is a constant satisfying

$$c = \epsilon \exp \left( -a - b \mu^T \theta^o - \frac{b^2 |\theta^o|^2}{2} \right) \tag{202}$$

$$= \exp \left( [a^* - a] + \mu^T [b^* \theta^* - b \theta^o] + \frac{b^{*2} |\theta^*|^2 - b^2 |\theta^o|^2}{2} \right). \tag{203}$$

We thus find Eq. (66). Note that the right hand side of the above equation depends implicitly on  $c$  through the value of  $\theta^o$ . To eliminate  $\theta^o$  from the equation for  $c$ , we let  $\Lambda$  denote the  $\theta^*$ -dependent part of Eq. (203), then substitute in Eq. (66):

$$\Lambda \equiv \mu^T [b^* \theta^* - b \theta^o] + \frac{b^{*2} |\theta^*|^2 - b^2 |\theta^o|^2}{2} \tag{204}$$

$$= \mu^T [b^* \theta^* (1 - c) - (c - 1) \mu] + \frac{b^{*2} |\theta^*|^2}{2} - \frac{|c b^* \theta^* + (c - 1) \mu|^2}{2} \tag{205}$$

$$= (1 - c) b^* \mu^T \theta^* + (1 - c) |\mu|^2 + \frac{(1 - c^2) b^{*2} |\theta^*|^2}{2} - \frac{(1 - c)^2 |\mu|^2}{2} - c(c - 1) b^* \mu^T \theta^*. \tag{206}$$

Using

$$(c - 1) - c(c - 1) = 1 - c^2, \quad \text{and} \quad (1 - c) - \frac{(1 - c)^2}{2} = \frac{1 - c^2}{2}, \tag{207}$$

we get

$$\Lambda = (1 - c^2) b^* \mu^T \theta^* + \frac{(1 - c^2) |\mu|^2}{2} + \frac{(1 - c^2) b^{*2} |\theta^*|^2}{2} \tag{208}$$

$$= \frac{1 - c^2}{2} |b^* \theta^* + \mu|^2. \tag{209}$$

We thus find the transcendental equation for  $c$ ,

$$c = \exp \left( [a^* - a] + \frac{1 - c^2}{2} |b^* \theta^* + \mu|^2 \right), \tag{210}$$

which is Eq. (67).

**Derivation of Eq. (70)**

From the expression for likelihood in Eq. (64), we find that the Hessian of likelihood is

$$H_{ij} = \left. \frac{\partial^2 L(\theta, a^*, b^*)}{\partial \theta_i \partial \theta_j} \right|_{\theta^*} \tag{211}$$

$$= [-b^{*2} \delta_{ij} - (b^* \mu_i + b^{*2} \theta_i^*)(b^* \mu_j + b^{*2} \theta_j^*)] \exp \left( a + b^* \theta^T \mu + \frac{b^{*2} |\theta|^2}{2} \right) \tag{212}$$

$$= -b^{*2} \epsilon \Lambda_{ij} \tag{213}$$

where

$$\Lambda_{ij} \equiv \delta_{ij} + (\mu_i + b^* \theta_i^*)(\mu_j + b^* \theta_j^*). \tag{214}$$

We note that in deriving Eq. (213) we used the expression for  $\epsilon$  in Eq. (58). The expression in Eq. (70) further makes use of the approximation  $N_1 \approx \epsilon N$ , which will hold in the  $\epsilon \rightarrow 0$  limit, and

$$\left. \frac{\partial^2 L(\theta, a^*, b^*)}{\partial \theta_i \partial \theta_j} \right|_{\theta^o} \approx \left. \frac{\partial^2 L(\theta, a^*, b^*)}{\partial \theta_i \partial \theta_j} \right|_{\theta^*}, \tag{215}$$

which will hold in the large  $N$  limit.

**Derivation of Eqs. (73) and (74)**

We derive Eq. (73) as follows. To ease notation a bit, we define  $p_M(R) = p(R|M)$ .

$$I[R; M] = \sum_{M=0,1} \int dR p(M, R) \log \frac{p_M(R)}{p(R)} \tag{216}$$

$$= p(M = 1) \int dR p_1(R) \log \frac{p_1(R)}{p(R)} + p(M = 0) \int dR p_0(R) \log \frac{p_0(R)}{p(R)} \tag{217}$$

$$= p(M = 1) \int dR p_1(R) \log \frac{p_1(R)}{p_0(R)} + p(M = 1) \int dR p_1(R) \log \frac{p_0(R)}{p(R)} + p(M = 0) \int dR p_0(R) \log \frac{p_0(R)}{p(R)} \tag{218}$$

$$= p(M = 1) \int dR p_1(R) \log \frac{p_1(R)}{p_0(R)} + \int dR p(R) \log \frac{p_0(R)}{p(R)}. \tag{219}$$

Because  $p(M = 1) = \epsilon + O(\epsilon^2)$ , the first term in Eq. (219) is the right hand side of Eq. (73) to lowest order in  $\epsilon$ . We now show that the second term is of order  $\epsilon^2$  and can therefore be ignored. Up to terms of order  $\epsilon^2$ ,

$$p(R) = (1 - \epsilon)p_0(R) + \epsilon p_1(R). \tag{220}$$

Rearranging this gives

$$p_0(R) = \frac{p(R) - \epsilon p_1(R)}{1 - \epsilon}. \tag{221}$$

Plugging this into the second term of Eq. (219) gives

$$\int dR p(R) \log \frac{p_0(R)}{p(R)} = \int dR p(R) \log \left[ \frac{1}{1 - \epsilon} \left( 1 - \epsilon \frac{p_1(R)}{p(R)} \right) \right] \tag{222}$$

$$= \int dR p(R) \log \left[ 1 + \epsilon \left( 1 - \frac{p_1(R)}{p(R)} \right) + O(\epsilon^2) \right] \tag{223}$$

$$= \epsilon \int dR p(R) \left( 1 - \frac{p_1(R)}{p(R)} \right) + O(\epsilon^2) \tag{224}$$

$$= O(\epsilon^2). \tag{225}$$

Equation (74) is derived as follows:

$$I(\theta) = \epsilon \left\langle \log \frac{p(R|M=1)}{p(R|M=0)} \right\rangle_{M=1} \tag{226}$$

$$= \epsilon \left\langle \frac{(R - \mu^T \theta)^2}{2|\theta|^2} - \frac{([R - \mu^T \theta] - b^* \theta^T \theta^*)^2}{2|\theta|^2} \right\rangle_{M=1} \tag{227}$$

$$= \frac{\epsilon}{2|\theta|^2} \left\langle 2[R - \mu^T \theta] b^* \theta^T \theta^* - (b^* \theta^T \theta^*)^2 \right\rangle_{M=1} \tag{228}$$

$$= \frac{\epsilon}{2|\theta|^2} \left( 2\langle [R]_{M=1} - \mu^T \theta \rangle b^* \theta^T \theta^* - (b^* \theta^T \theta^*)^2 \right) \tag{229}$$

$$= \frac{\epsilon}{2|\theta|^2} \left( 2[b^* \theta^T \theta^*] b^* \theta^T \theta^* - (b^* \theta^T \theta^*)^2 \right) \tag{230}$$

$$= \frac{\epsilon b^{*2}}{2} \frac{(\theta^T \theta^*)^2}{|\theta|^2}. \tag{231}$$

**Derivation of Eq. (77)**

To derive Eq. (77), we set

$$\theta = \theta^* + \delta\theta_{\parallel} + \delta\theta_{\perp} \tag{232}$$

where  $\delta\theta_{\parallel}$  is the deviation of  $\theta$  from  $\theta^*$  in the direction of  $\theta^*$ , and  $\delta\theta_{\perp}$  is the deviation perpendicular to  $\theta^*$ . This gives

$$\frac{(\theta^T \theta^*)^2}{|\theta|^2} = \frac{(|\theta^*|^2 + \delta\theta_{\parallel}^T \theta^*)^2}{|\theta^*|^2 + 2\delta\theta_{\parallel}^T \theta^* + |\delta\theta_{\parallel}|^2 + |\delta\theta_{\perp}|^2} \tag{233}$$

$$= |\theta^*|^2 \frac{|\theta^*|^2 + 2\delta\theta_{\parallel}^T \theta^* + |\delta\theta_{\parallel}|^2}{|\theta^*|^2 + 2\delta\theta_{\parallel}^T \theta^* + |\delta\theta_{\parallel}|^2 + |\delta\theta_{\perp}|^2} \tag{234}$$

$$= |\theta^*|^2 \left( 1 - \frac{|\delta\theta_{\perp}|^2}{|\theta^*|^2} + \dots \right) \tag{235}$$

$$= |\theta^*|^2 - |\delta\theta_{\perp}|^2 + \dots \tag{236}$$

The result in Eq. (77) readily follows by substituting this into the formula for mutual information in Eq. (74), then approximating the Hessian of mutual information at  $\theta^o$  by the Hessian at  $\theta^*$ .

**Derivation of Eq. 80**

Here we show how to evaluate the equation, Eq. (79), for the noise-averaged likelihood  $e^{NL_{na}(\theta)}$ . First, interchange the order of integration and define  $a' = a + b\theta^T \mu$ . This gives,

$$e^{NL_{na}(\theta)} = C \int_{-\infty}^{\infty} db \int_{-\infty}^{\infty} da' \exp \left[ N_1 \left[ a' + bb^* \theta^T \theta^* \right] - N \exp \left( a' + \frac{b^2 |\theta|^2}{2} \right) \right]. \tag{237}$$

Next, define  $M = N \exp \left( \frac{b^2 |\theta|^2}{2} \right)$ ,  $u = Me^{a'}$ , and so  $e^{a'} = u/M$ ,  $e^{a'} da' = du/M$ . This gives

$$e^{NL_{na}(\theta)} = C \int_{-\infty}^{\infty} db e^{N_1 bb^* \theta^T \theta^*} \int_{-\infty}^{\infty} (e^{a'} da') (e^{a'})^{N_1 - 1} \exp[-Me^{a'}] \tag{238}$$

$$= C \int_{-\infty}^{\infty} db e^{N_1 bb^* \theta^T \theta^*} M^{-N_1} \int_0^{\infty} du u^{N_1 - 1} \exp[-u] \tag{239}$$

$$= C \Gamma(N_1) \int_{-\infty}^{\infty} db e^{N_1 bb^* \theta^T \theta^*} M^{-N_1} \tag{240}$$

$$= C \Gamma(N_1) \int_{-\infty}^{\infty} db \exp \left[ N_1 bb^* \theta^T \theta^* - N_1 \frac{b^2 |\theta|^2}{2} \right] \tag{241}$$

$$= C \Gamma(N_1) \int_{-\infty}^{\infty} db \exp \left[ \frac{N_1 |\theta|^2}{2} \left( 2bb^* \frac{\theta^T \theta^*}{|\theta|^2} - b^2 \right) \right] \tag{242}$$

$$= C \Gamma(N_1) \int_{-\infty}^{\infty} db \exp \left[ \frac{N_1 |\theta|^2}{2} \left( \frac{b^{*2} (\theta^T \theta^*)^2}{|\theta|^4} - \left[ b - \frac{b^* \theta^T \theta^*}{|\theta|^2} \right]^2 \right) \right] \tag{243}$$

$$= C \Gamma(N_1) \exp \left( \frac{N_1 b^{*2} (\theta^T \theta^*)^2}{2 |\theta|^2} \right) \int_{-\infty}^{\infty} db \exp \left( \frac{N_1 |\theta|^2}{2} \left[ b - \frac{b^* \theta^T \theta^*}{|\theta|^2} \right]^2 \right) \tag{244}$$

$$= C \Gamma(N_1) \sqrt{\frac{2\pi}{N_1 |\theta|^2}} \exp \left( \frac{N_1 b^{*2} (\theta^T \theta^*)^2}{2 |\theta|^2} \right), \tag{245}$$

which is Eq. (80).

**References**

1. Melnikov, A., Murugan, A., Zhang, X., Tesileanu, T., Wang, L., Rogov, P., Feizi, S., Gnirke, A., Callan, C.G., Kinney, J.B., Kellis, M., Lander, E.S., Mikkelsen, T.S.: Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**(3), 271–277 (2012)
2. Mukherjee, S., Berger, M., Jona, G., Wang, X., Muzzey, D., Snyder, M., Young, R., Bulyk, M.: Rapid analysis of the DNA-binding specificities of transcription factors with DNA microarrays. *Nat. Genet.* **36**(12), 1331–1339 (2004)

3. Berger, M., Philippakis, A., Qureshi, A., He, F., Estep, P., Bulyk, M.: Compact, universal DNA microarrays to comprehensively determine transcription-factor binding site specificities. *Nat. Biotechnol.* **24**(11), 1429–1435 (2006)
4. Meng, X., Brodsky, M.H., Wolfe, S.A.: A bacterial one-hybrid system for determining the DNA-binding specificity of transcription factors. *Nat. Biotechnol.* **23**(8), 988–994 (2005)
5. Maerkl, S., Quake, S.: A systems approach to measuring the binding energy landscapes of transcription factors. *Science* **315**(5809), 233–237 (2007)
6. Zykovich, A., Korf, I., Segal, D.J.: Bind-n-Seq: high-throughput analysis of in vitro protein-DNA interactions using massively parallel sequencing. *Nucleic Acids Res.* **37**(22), e151 (2009)
7. Zhao, Y., Granas, D., Stormo, G.D.: Inferring binding energies from selected binding sites. *PLoS Comput. Biol.* **5**(12), e1000590 (2009)
8. Jolma, A., Kivioja, T., Toivonen, J., Cheng, L., Wei, G., Enge, M., Taipale, M., Vaquerizas, J.M., Yan, J., Sillanpää, M.J., Bonke, M., Palin, K., Talukder, S., Hughes, T.R., Luscombe, N.M., Ukkonen, E., Taipale, J.: Multiplexed massively parallel SELEX for characterization of human transcription factor binding specificities. *Genome Res.* **20**(6), 861–873 (2010)
9. Wong, D., Teixeira, A., Oikonomopoulos, S., Humburg, P., Lone, I.N., Saliba, D., Siggers, T., Bulyk, M., Angelov, D., Dimitrov, S., Udalova, I.A., Ragoussis, J.: Extensive characterization of NF- $\kappa$ B binding uncovers non-canonical motifs and advances the interpretation of genetic functional traits. *Genome Biol.* **12**(7), R70 (2011)
10. Slattery, M., Riley, T., Liu, P., Abe, N., Gomez-Alcala, P., Dror, I., Zhou, T., Rohs, R., Honig, B., Bussemaker, H.J., Mann, R.S.: Cofactor binding evokes latent differences in DNA binding specificity between Hox proteins. *Cell* **147**(6), 1270–1282 (2011)
11. Patwardhan, R.P., Lee, C., Litvin, O., Young, D.L., Pe'er, D., Shendure, J.: High-resolution analysis of DNA regulatory elements by synthetic saturation mutagenesis. *Nat. Biotechnol.* **27**(12), 1173–1175 (2009)
12. Kinney, J.B., Murugan, A., Callan, C.G., Cox, E.C.: Using deep sequencing to characterize the biophysical mechanism of a transcriptional regulatory sequence. *Proc. Natl. Acad. Sci. USA* **107**(20), 9158–9163 (2010)
13. Patwardhan, R.P., Hiatt, J.B., Witten, D.M., Kim, M.J., Smith, R.P., May, D., Lee, C., Andrie, J.M., Lee, S.I., Cooper, G.M., Ahituv, N., Pennacchio, L.A., Shendure, J.: Massively parallel functional dissection of mammalian enhancers in vivo. *Nat. Biotechnol.* **30**(3), 265–270 (2012)
14. Sharon, E., Kalma, Y., Sharp, A., Raveh-Sadka, T., Levo, M., Zeevi, D., Keren, L., Yakhini, Z., Weinberger, A., Segal, E.: Inferring gene regulatory logic from high-throughput measurements of thousands of systematically designed promoters. *Nat. Biotechnol.* **30**(6), 521–530 (2012)
15. Kwasnieski, J.C., Mogno, I., Myers, C.A., Corbo, J.C., Cohen, B.A.: Complex effects of nucleotide variants in a mammalian cis-regulatory element. *Proc. Natl. Acad. Sci. USA* **109**(47), 19498–19503 (2012)
16. Fowler, D.M., Araya, C.L., Fleishman, S.J., Kellogg, E.H., Stephany, J.J., Baker, D., Fields, S.: High-resolution mapping of protein sequence-function relationships. *Nat. Methods* **7**(9), 741–746 (2010)
17. Hietpas, R.T., Jensen, J.D., Bolon, D.N.A.: Experimental illumination of a fitness landscape. *Proc. Natl. Acad. Sci. USA* **108**(19), 7896–7901 (2011)
18. Adkar, B.V., Tripathi, A., Sahoo, A., Bajaj, K., Goswami, D., Chakrabarti, P., Swarnkar, M.K., Gokhale, R.S., Varadarajan, R.: Protein model discrimination using mutational sensitivity derived from deep sequencing. *Structure* **20**(2), 371–381 (2012)
19. Whitehead, T.A., Chevalier, A., Song, Y., Dreyfus, C., Fleishman, S.J., De Mattos, C., Myers, C.A., Kamisetty, H., Blair, P., Wilson, I.A., Baker, D.: Optimization of affinity, specificity and function of designed influenza inhibitors using deep sequencing. *Nat. Biotechnol.* **30**(6), 543–548 (2012)
20. Schlinkmann, K.M., Honegger, A., Türeci, E., Robison, K.E., Lipovšek, D., Plückthun, A.: Critical features for biosynthesis, stability, and functionality of a G protein-coupled receptor uncovered by all-versus-all mutations. *Proc. Natl. Acad. Sci. USA* **109**(25), 9810–9815 (2012)
21. Holmqvist, E., Reimegård, J., Wagner, E.G.H.: Massive functional mapping of a 5'-UTR by saturation mutagenesis, phenotypic sorting and deep sequencing. *Nucleic Acids Res.* **41**(12), e122 (2013)
22. Peterman, N., Lavi-Itzkovitz, A., Levine, E.: Large-scale mapping of sequence-function relations in small regulatory RNAs reveals plasticity and modularity. *Nucleic Acids Res.* **42**(19), 12177–12188 (2014)
23. Oikonomou, P., Goodarzi, H., Tavazoie, S.: Systematic identification of regulatory elements in conserved 3' UTRs of human transcripts. *Cell Rep.* **7**(1), 281–292 (2014)
24. Noderer, W.L., Flockhart, R.J., Bhaduri, A., Diaz de Arce, A.J., Zhang, J., Khavari, P.A., Wang, C.L.: Quantitative analysis of mammalian translation initiation sites by FACS-seq. *Mol. Syst. Biol.* **10**(8), 748 (2014)



25. Liachko, I., Youngblood, R.A., Keich, U., Dunham, M.J.: High-resolution mapping, characterization, and optimization of autonomously replicating sequences in yeast. *Genome Res.* **23**(4), 698–704 (2013)
26. Thyme, S.B., Song, Y., Brunette, T.J., Szeto, M.D., Kusak, L., Bradley, P., Baker, D.: Massively parallel determination and modeling of endonuclease substrate specificity. *Nucleic Acids Res.* **42**(22), 13839–13852 (2014)
27. Kinney, J.B., Tkacik, G., Callan, C.G.: Precise physical models of protein-DNA interaction from high-throughput data. *Proc. Natl. Acad. Sci. USA* **104**(2), 501–506 (2007)
28. Kinney, J.B., Atwal, G.S.: Parametric inference in the large data limit using maximally informative models. *Neural Comput.* **26**(4), 637–653 (2014)
29. Cover, T., Thomas, J.: *Elements of Information Theory*, 1st edn. Wiley, New York (1991)
30. Paninski, L.: Convergence properties of three spike-triggered analysis techniques. *Network-Comput. Neural* **14**(3), 437–464 (2003)
31. Sharpee, T., Rust, N., Bialek, W.: Analyzing neural responses to natural signals: maximally informative dimensions. *Neural Comput.* **16**(2), 223–250 (2004)
32. Sharpee, T., Sugihara, H., Kurgansky, A., Rebrik, S., Stryker, M., Miller, K.: Adaptive filtering enhances information transmission in visual cortex. *Nature* **439**(7079), 936–942 (2006)
33. Kouh, M., Sharpee, T.O.: Estimating linear-nonlinear models using Rényi divergences. *Network-Comput. Neural* **20**(2), 49–68 (2009)
34. Rajan, K., Marre, O., Tkacik, G.: Learning quadratic receptive fields from neural responses to natural signals: information theoretic and likelihood methods. *Neural Comput.* **25**(7), 1661–1692 (2013)
35. Jolma, A., Yan, J., Whittington, T., Toivonen, J., Nitta, K.R., Rastas, P., Morgunova, E., Enge, M., Taipale, M., Wei, G., Palin, K., Vaquerizas, J.M., Vincentelli, R., Luscombe, N.M., Hughes, T.R., Lemaire, P., Ukkonen, E., Kivioja, T., Taipale, J.: DNA-binding specificities of human transcription factors. *Cell* **152**(1), 327–339 (2013)
36. Oliphant, A., Brandl, C., Struhl, K.: Defining the sequence specificity of DNA-binding proteins by selecting binding sites from random-sequence oligonucleotides: analysis of yeast GCN4 protein. *Mol. Cell. Biol.* **9**(7), 2944–2949 (1989)
37. Tuerk, C., Gold, L.: Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage T4 DNA polymerase. *Science* **249**(4968), 505–510 (1990)
38. Ellington, A.D., Szostak, J.W.: In vitro selection of RNA molecules that bind specific ligands. *Nature* **346**(6287), 818–822 (1990)
39. Blackwell, T.K., Weintraub, H.: Differences and similarities in DNA-binding preferences of MyoD and E2A protein complexes revealed by binding site selection. *Science* **250**(4984), 1104–1110 (1990)
40. Wright, W., Binder, M., Funk, W.: Cyclic amplification and selection of targets (CASTing) for the myogenin consensus binding site. *Mol. Cell. Biol.* **11**(8), 4104–4110 (1991)
41. Herzenberg, L., Sweet, R., Herzenberg, L.: Fluorescence-activated cell sorting. *Sci. Am.* **234**(3), 108–117 (1976)
42. Fowler, D.M., Fields, S.: Deep mutational scanning: a new style of protein science. *Nat. Methods* **11**(8), 801–807 (2014)
43. Kinney, J.B., Atwal, G.S.: Equitability, mutual information, and the maximal information coefficient. *Proc. Natl. Acad. Sci. USA* **111**(9), 3354–3359 (2014)
44. Khan, S., Bandyopadhyay, S., Ganguly, A., Saigal, S., Erickson III, D., Protopopescu, V., Ostrouchov, G.: Relative performance of mutual information estimation methods for quantifying the dependence among short and noisy data. *Phys. Rev. E* **76**(2), 026209 (2007)
45. Bialek, W., Callan, C., Strong, S.: Field theories for learning probability distributions. *Phys. Rev. Lett.* **77**(23), 4693–4697 (1996)
46. Kinney, J.B.: Estimation of probability densities using scale-free field theories. *Phys. Rev. E* **90**(1), 011301(R) (2014)
47. Kinney, J.B.: Unification of field theory and maximum entropy methods for learning probability densities. *Phys. Rev. E* **92**(3), 032107 (2015)