

# Predictive information: From definition to applications to biological systems

Ilya Nemenman

KITP, UCSB

Thanks to: William Bialek, Naftali Tishby

physics/0007070

physics/0103076

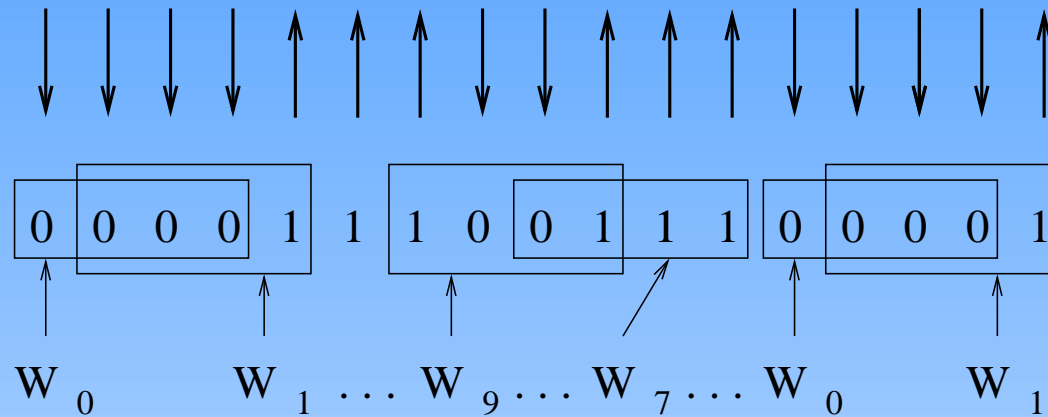
q-bio/0402029



# Outline

- A curious observation.
- Quantifying predictability and complexity.
- Predictability and optimization in sensory information processing.
- Learning and predictive information.
- Testing models used by animals.
- Bonus material.

# Entropy of words in a spin chain



$$W_0 = 0 \ 0 \ 0 \ 0$$

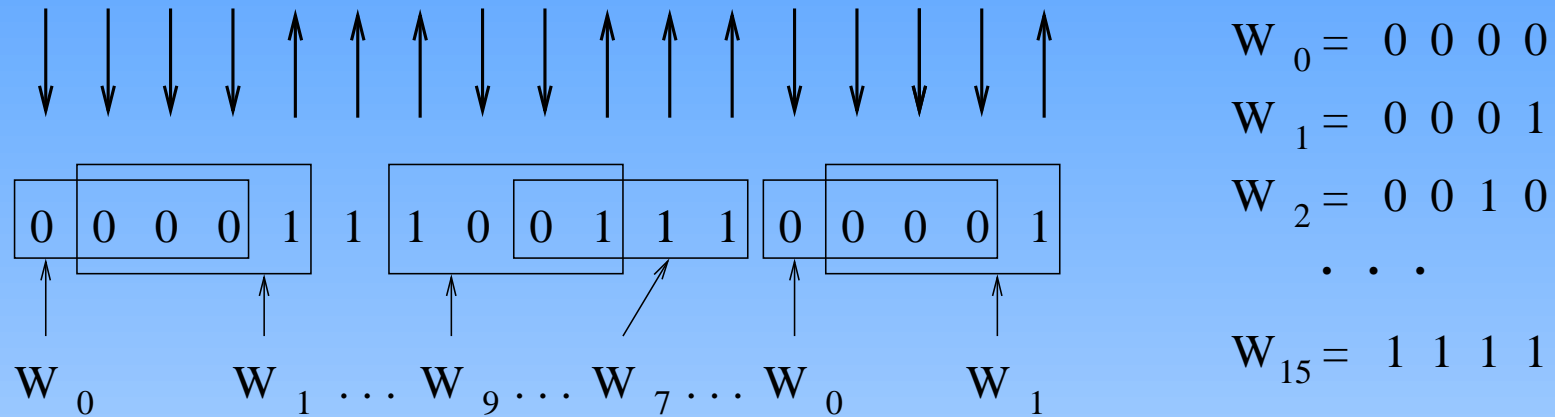
$$W_1 = 0 \ 0 \ 0 \ 1$$

$$W_2 = 0 \ 0 \ 1 \ 0$$

...

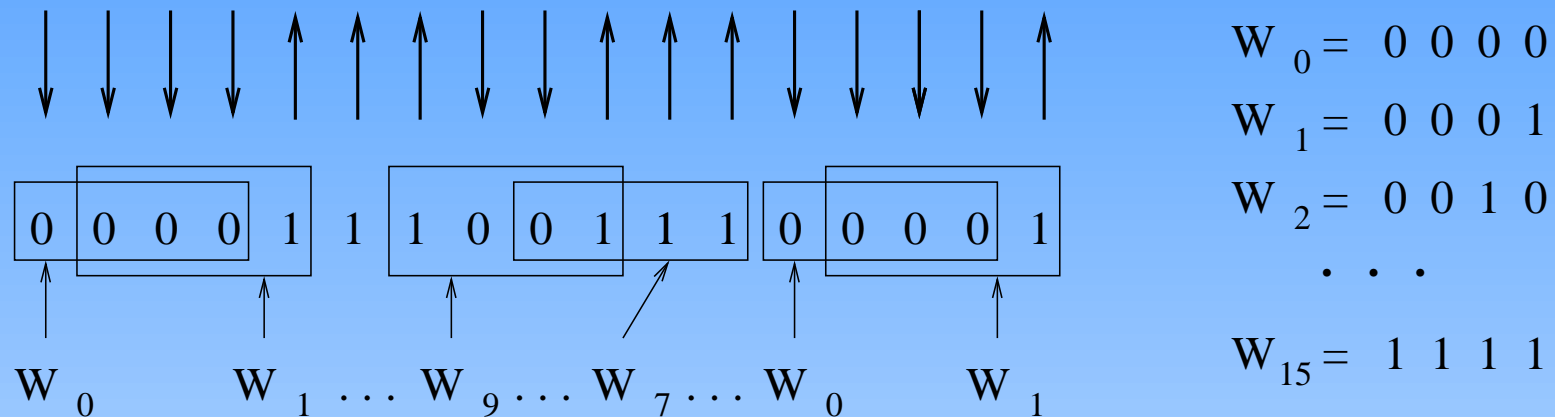
$$W_{15} = 1 \ 1 \ 1 \ 1$$

# Entropy of words in a spin chain



$$S(N) = - \sum_{k=0}^{2^N-1} P_N(W_k) \log_2 P_N(W_k)$$

# Entropy of words in a spin chain

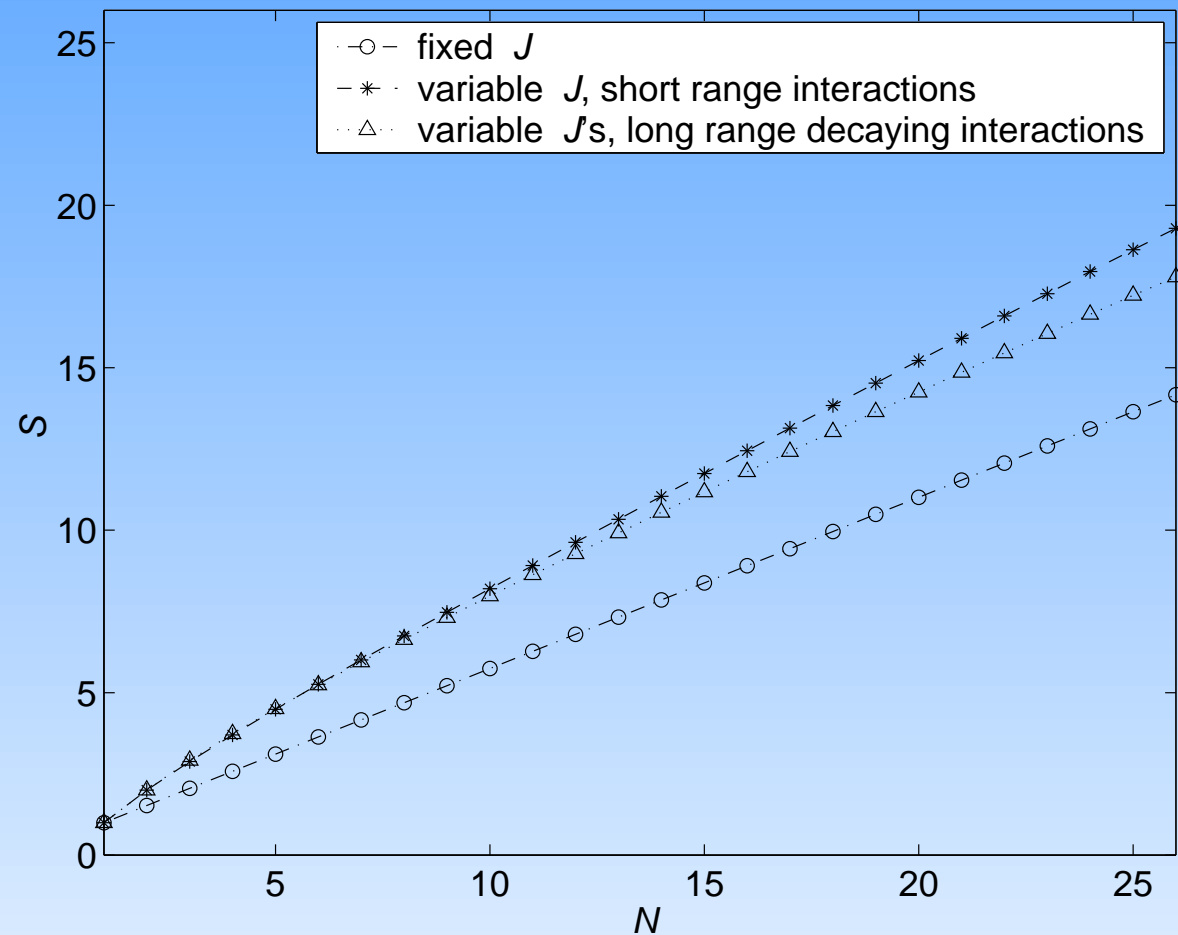


$$S(N) = - \sum_{k=0}^{2^N-1} P_N(W_k) \log_2 P_N(W_k)$$

For this chain,  $P(W_0) = P(W_1) = P(W_3) = P(W_7) = P(W_{12}) = P(W_{14}) = 2$ ,  $P(W_8) = P(W_9) = 1$ , and all other frequencies (probabilities) are zero. Thus,  $S(4) \approx 2.95$  bits.

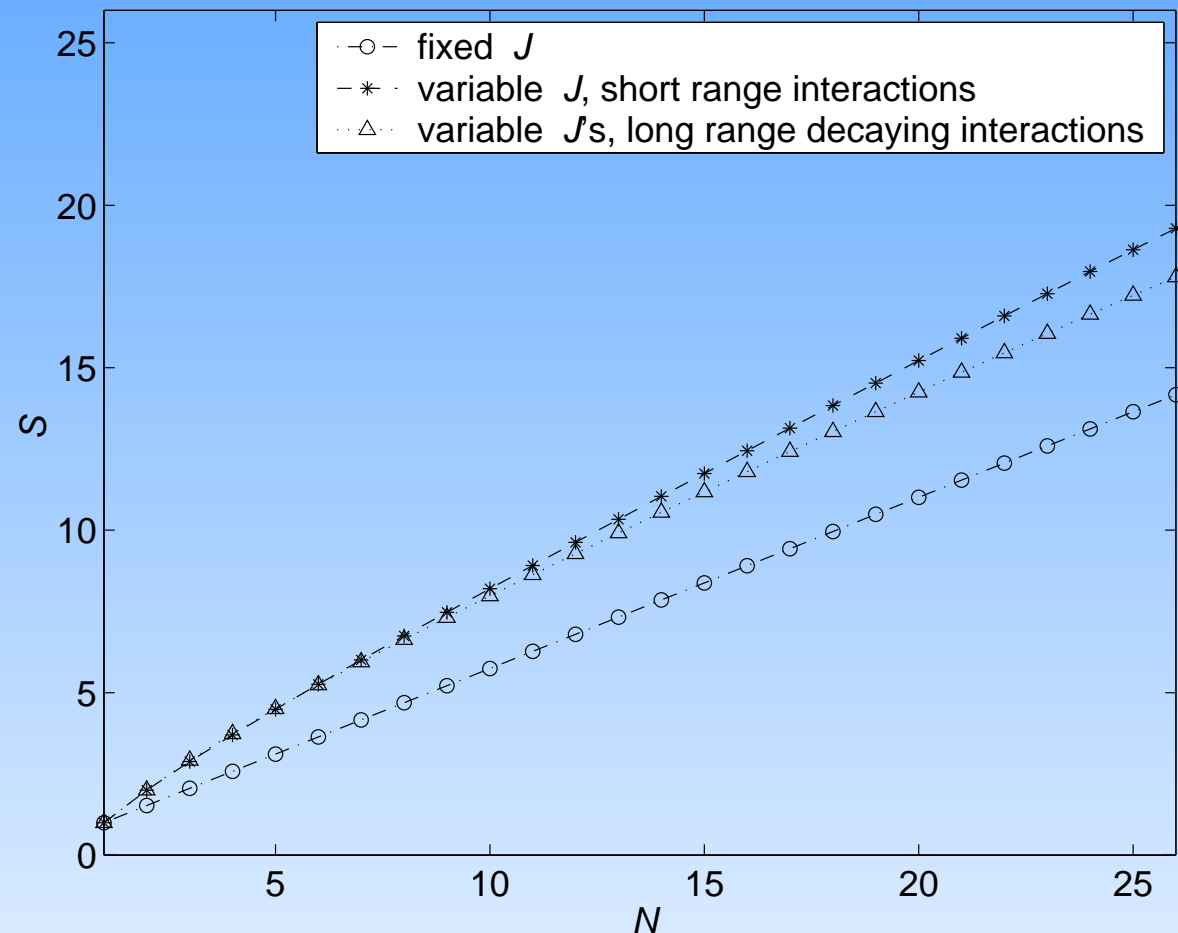
# Entropy of 3 generated chains

- $J_{ij} = \delta_{i,j+1}$
  - $J_{ij} = J_0 \delta_{i,j+1}$ ,  $J_0$  is taken at random from  $\mathcal{N}(0, 1)$  every 400000 spins
  - $J_{ij}$  is taken at random from  $\mathcal{N}(0, \frac{1}{i-j})$  every 400000 spins
- $1 \cdot 10^9$  spins total.



# Entropy of 3 generated chains

- $J_{ij} = \delta_{i,j+1}$
  - $J_{ij} = J_0 \delta_{i,j+1}$ ,  $J_0$  is taken at random from  $\mathcal{N}(0, 1)$  every 400000 spins
  - $J_{ij}$  is taken at random from  $\mathcal{N}(0, \frac{1}{i-j})$  every 400000 spins
- $1 \cdot 10^9$  spins total.

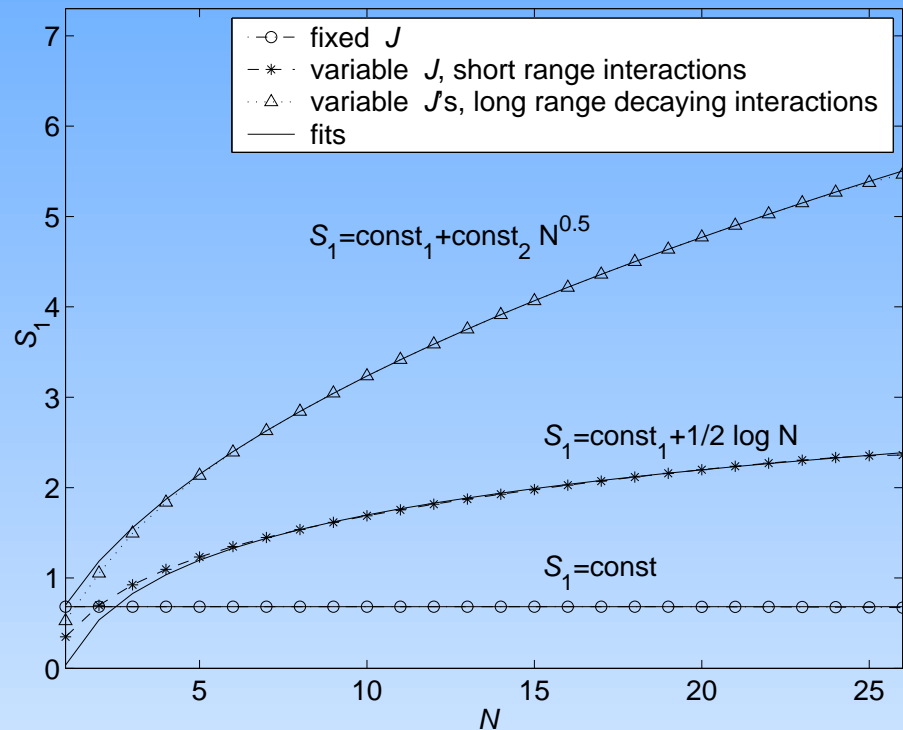


Entropy is extensive!

It shows no distinction between the cases.

# Subextensive component of the entropy

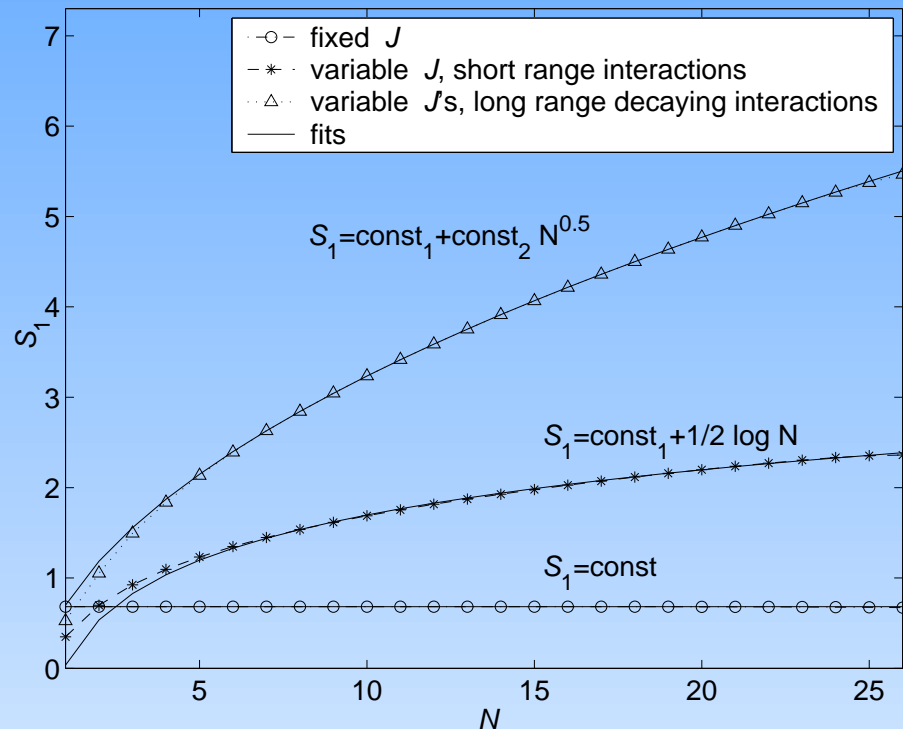
... shows a qualitative distinction between the cases!





# Subextensive component of the entropy

... shows a qualitative distinction between the cases!



Other examples:

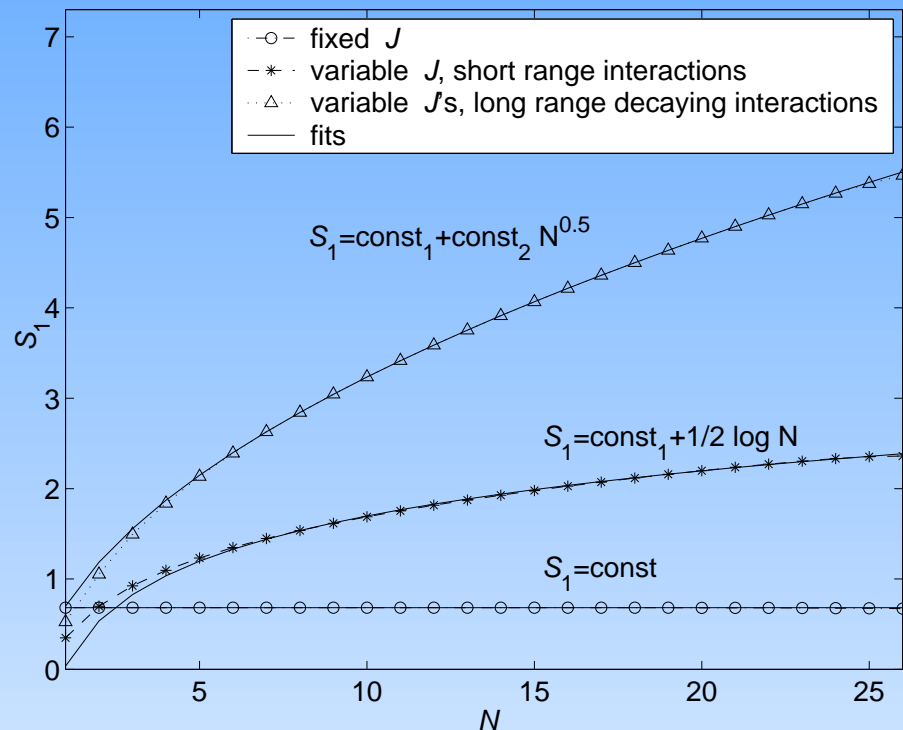
**const** periodic, fully random, chaotic sequences (finite correlation length)

**log** systems at phase transitions, or at the onset of chaos (divergent correlation length)

**power** texts, DNA sequences, (likely) some exotic transitions, (many divergent correlation lengths)

# Subextensive component of the entropy

... shows a qualitative distinction between the cases!



Other examples:

**const** periodic, fully random, chaotic sequences (finite correlation length)

**log** systems at phase transitions, or at the onset of chaos (divergent correlation length)

**power** texts, DNA sequences, (likely) some exotic transitions, (many divergent correlation lengths)

- Entropy density or channel capacity do not distinguish these cases.
- Theory of phase transitions may not distinguish between the last two cases.
- Complexity of underlying dynamics intuitively increases from **const** to **power**.

# Objectives

- unified description of complexity and learning
- make distinction between useful and unusable data
- do this using physical quantities
- understand models used by organisms to represent the world
- understand biological designs by means of optimization principles

## Solution – predictability

- we learn (estimate parameters, extrapolate, classify, ...) to *generalize* and *predict* from training examples; estimation of parameters is only an intermediate step

## Solution – predictability

- we learn (estimate parameters, extrapolate, classify, ...) to *generalize* and *predict* from training examples; estimation of parameters is only an intermediate step
- nonpredictive features in any signal are useless since we observe *now* and react in the *future*

## Solution – predictability

- we learn (estimate parameters, extrapolate, classify, ...) to *generalize* and *predict* from training examples; estimation of parameters is only an intermediate step
- nonpredictive features in any signal are useless since we observe *now* and react in the *future*
- high predictability sources (more details to predict, not easier predictions) are generated by more complex sources (in particular, regular and random sources have low complexity)

## Solution – predictability

- we learn (estimate parameters, extrapolate, classify, ...) to *generalize* and *predict* from training examples; estimation of parameters is only an intermediate step
- nonpredictive features in any signal are useless since we observe *now* and react in the *future*
- high predictability sources (more details to predict, not easier predictions) are generated by more complex sources (in particular, regular and random sources have low complexity)
- measuring organisms' learning and prediction performance for signals of different complexity may reveal the underlying models

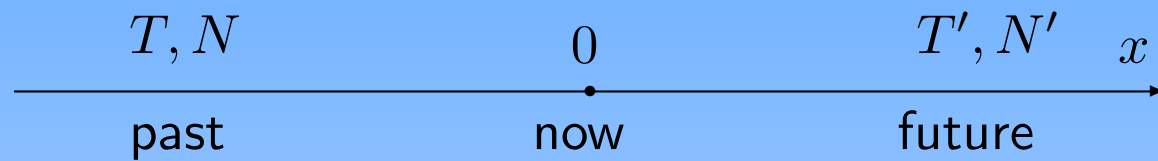
## Solution – predictability

- we learn (estimate parameters, extrapolate, classify, ...) to *generalize* and *predict* from training examples; estimation of parameters is only an intermediate step
- nonpredictive features in any signal are useless since we observe *now* and react in the *future*
- high predictability sources (more details to predict, not easier predictions) are generated by more complex sources (in particular, regular and random sources have low complexity)
- measuring organisms' learning and prediction performance for signals of different complexity may reveal the underlying models
- optimizing predictive information may be the design principle



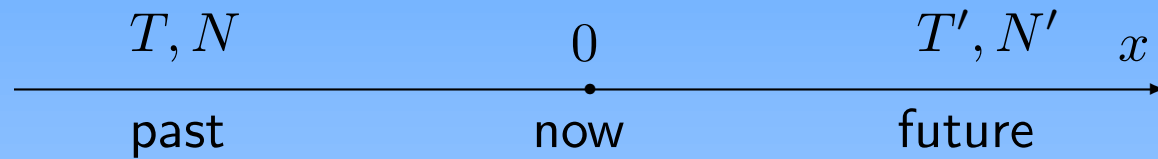
# Quantifying predictability

Information theory: non-metric, universal way to quantify learning



# Quantifying predictability

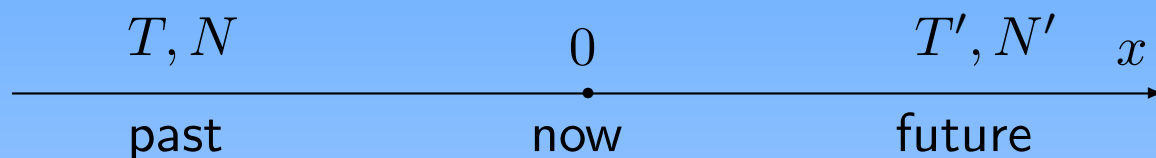
Information theory: non-metric, universal way to quantify learning



$$\begin{aligned}\mathcal{I}_{\text{pred}}(T, T') &= \left\langle \log_2 \left[ \frac{P(x_{\text{future}} | x_{\text{past}})}{P(x_{\text{future}})} \right] \right\rangle \\ &= S(T) + S(T') - S(T + T')\end{aligned}$$

# Quantifying predictability

Information theory: non-metric, universal way to quantify learning



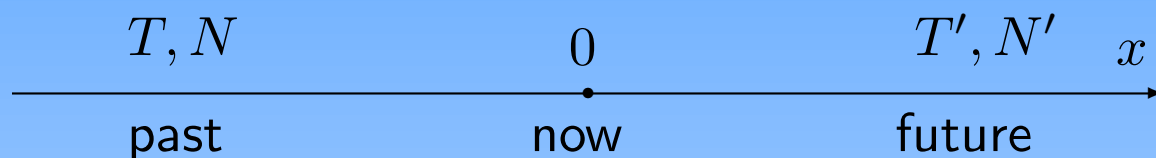
$$\begin{aligned}
 \mathcal{I}_{\text{pred}}(T, T') &= \left\langle \log_2 \left[ \frac{P(x_{\text{future}} | x_{\text{past}})}{P(x_{\text{future}})} \right] \right\rangle \\
 &= S(T) + S(T') - S(T + T') \\
 S(T) &= \mathcal{S}_0 \cdot T + S_1(T)
 \end{aligned}$$

Extensive component cancels in predictive information.

Predictability is a deviation from extensivity!

# Quantifying predictability

Information theory: non-metric, universal way to quantify learning



$$\begin{aligned}
 \mathcal{I}_{\text{pred}}(T, T') &= \left\langle \log_2 \left[ \frac{P(x_{\text{future}} | x_{\text{past}})}{P(x_{\text{future}})} \right] \right\rangle \\
 &= S(T) + S(T') - S(T + T') \\
 S(T) &= \mathcal{S}_0 \cdot T + S_1(T)
 \end{aligned}$$

Extensive component cancels in predictive information.

Predictability is a deviation from extensivity!

$$I_{\text{pred}}(T) \equiv \mathcal{I}_{\text{pred}}(T, \infty) = S_1(T)$$

## Properties of $I_{\text{pred}}(T)$

- $I_{\text{pred}}(T)$  is information, so  $I_{\text{pred}}(T) \geq 0$

## Properties of $I_{\text{pred}}(T)$

- $I_{\text{pred}}(T)$  is information, so  $I_{\text{pred}}(T) \geq 0$
- $I_{\text{pred}}(T)$  is subextensive,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{T} = 0$

## Properties of $I_{\text{pred}}(T)$

- $I_{\text{pred}}(T)$  is information, so  $I_{\text{pred}}(T) \geq 0$
- $I_{\text{pred}}(T)$  is subextensive,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{T} = 0$
- diminishing returns,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{S(T)} = 0$

## Properties of $I_{\text{pred}}(T)$

- $I_{\text{pred}}(T)$  is information, so  $I_{\text{pred}}(T) \geq 0$
- $I_{\text{pred}}(T)$  is subextensive,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{T} = 0$
- diminishing returns,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{S(T)} = 0$
- prediction and postdiction are symmetric



## Properties of $I_{\text{pred}}(T)$

- $I_{\text{pred}}(T)$  is information, so  $I_{\text{pred}}(T) \geq 0$
- $I_{\text{pred}}(T)$  is subextensive,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{T} = 0$
- diminishing returns,  $\lim_{T \rightarrow \infty} \frac{I_{\text{pred}}(T)}{S(T)} = 0$
- prediction and postdiction are symmetric
- it relates to and generalizes many relevant quantities
  - learning: universal learning curves
  - complexity: complexity measures
  - coding: model coding length

## How can $I_{\text{pred}}$ behave?

$\lim_{N \rightarrow \infty} I_{\text{pred}} = \text{const}$  no long-range structure

- simply predictable (periodic, constant, etc.) processes
- fully stochastic (Markov) processes

## How can $I_{\text{pred}}$ behave?

$\lim_{N \rightarrow \infty} I_{\text{pred}} = \text{const}$  no long-range structure

- simply predictable (periodic, constant, etc.) processes
- fully stochastic (Markov) processes

$\lim_{N \rightarrow \infty} I_{\text{pred}} = \text{const} \times \log_2 N$  precise learning of a fixed set of parameters

- learning finite-parameter densities
- well known as  $I(N, \text{parameters}) = I_{\text{pred}}(N)$
- physical system at criticality
- (possibly) nonextensive statistics systems

## How can $I_{\text{pred}}$ behave?

$\lim_{N \rightarrow \infty} I_{\text{pred}} = \text{const}$  no long-range structure

- simply predictable (periodic, constant, etc.) processes
- fully stochastic (Markov) processes

$\lim_{N \rightarrow \infty} I_{\text{pred}} = \text{const} \times \log_2 N$  precise learning of a fixed set of parameters

- learning finite-parameter densities
- well known as  $I(N, \text{parameters}) = I_{\text{pred}}(N)$
- physical system at criticality
- (possibly) nonextensive statistics systems

$\lim_{N \rightarrow \infty} I_{\text{pred}} = \text{const} \times N^\xi$  learning more features as  $N$  grows

- learning continuous densities
- language
- some critical phenomena (wetting transitions)
- not well studied

## Which complexity do we want to define?

- complexity of dynamics that generates a time series (not computational or descriptive complexity); thus it must be zero for totally random and for easily predictable processes
- usable for Occam-style punishment in statistical inference
- expressible in conventional physical terms
- must be attached to an ensemble, not a single realization

## Complexity measure

- some kind of entropy (we proclaim Shannon's postulates: monotonicity, continuity, additivity)

## Complexity measure

- some kind of entropy (we proclaim Shannon's postulates: monotonicity, continuity, additivity)
- invariant under invertible temporally local transformations ( $x_k \rightarrow x_k + \xi x_{k-1}$ : measuring device with inertia, article with misprints, same book in different languages – same universality class)

$$\log P_1(x) = \log P_2(x) + \text{loc. oper.} \Rightarrow C[P_1(x)] = C[P_2(x)]$$

This may present a problem in higher dimensions.

# Complexity measure

- some kind of entropy (we proclaim Shannon's postulates: monotonicity, continuity, additivity)
- invariant under invertible temporally local transformations ( $x_k \rightarrow x_k + \xi x_{k-1}$ : measuring device with inertia, article with misprints, same book in different languages – same universality class)

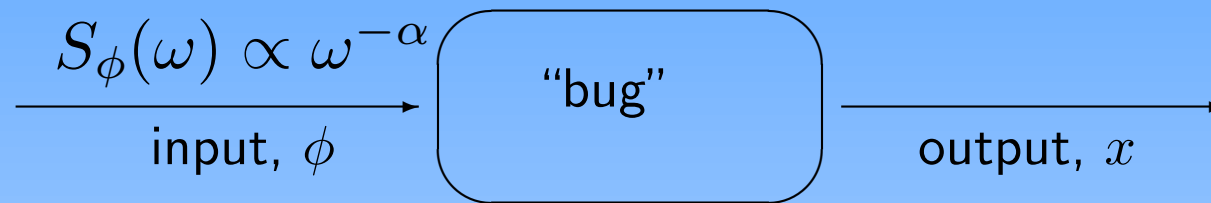
$$\log P_1(x) = \log P_2(x) + \text{loc. oper.} \Rightarrow C[P_1(x)] = C[P_2(x)]$$

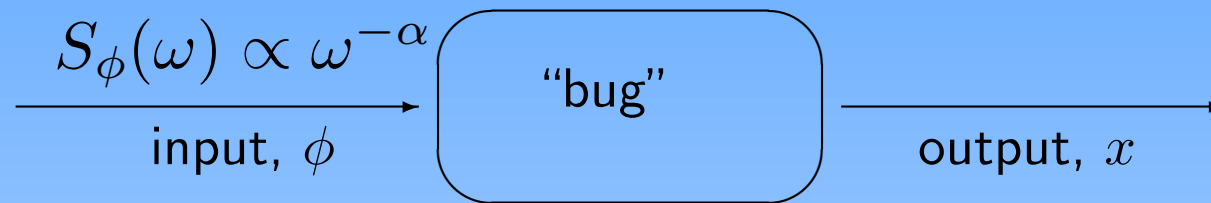
This may present a problem in higher dimensions.

The divergent subextensive term measures complexity uniquely!



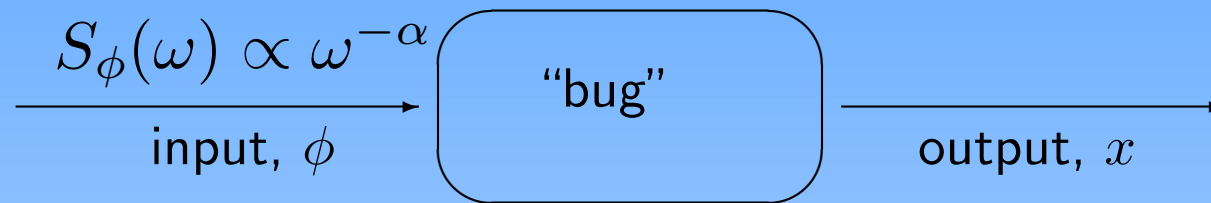
# $I_{\text{pred}}$ optimization in biology



$I_{\text{pred}}$  optimization in biology

$$\tau \frac{dx}{dt} = -x + \phi(t) + \eta(t), \quad \langle \eta(t) \eta(0) \rangle = 1/I_0 \delta(t)$$

# $I_{\text{pred}}$ optimization in biology

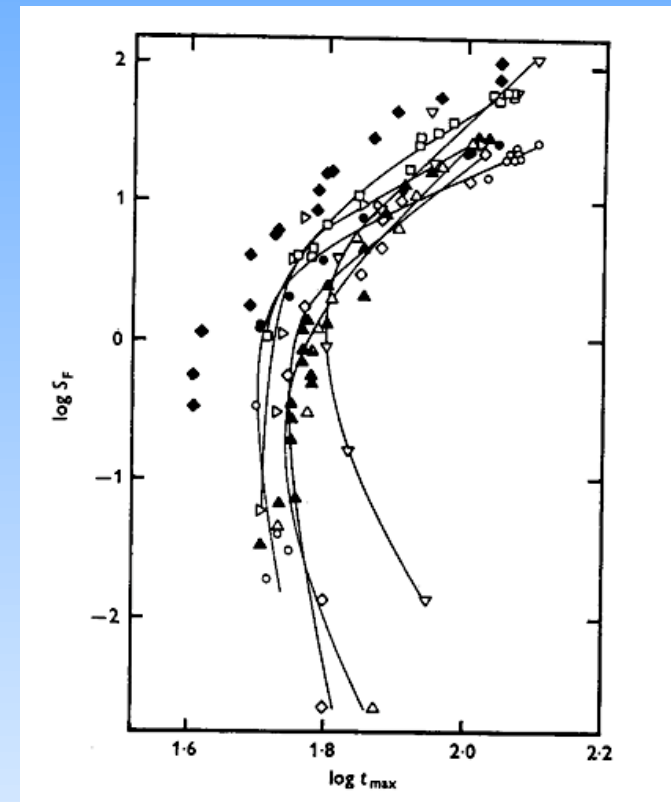


$$\tau \frac{dx}{dt} = -x + \phi(t) + \eta(t), \quad \langle \eta(t)\eta(0) \rangle = 1/I_0 \delta(t)$$

$$\mathcal{I}([\phi], [x]) = \lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T/2}^{T/2} \frac{d\omega}{2\pi} \log \left( 1 + \frac{S_\phi(\omega)}{1/I_0} \right)$$

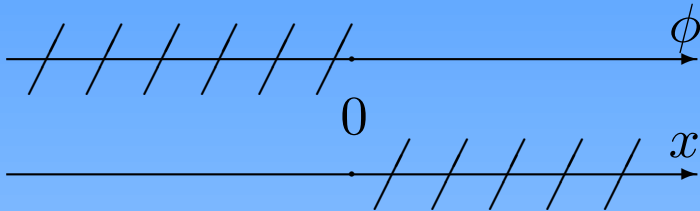
Maximization w.r.t.  $\tau$  is meaningless.

# $I_{\text{pred}}$ extraction and maximization

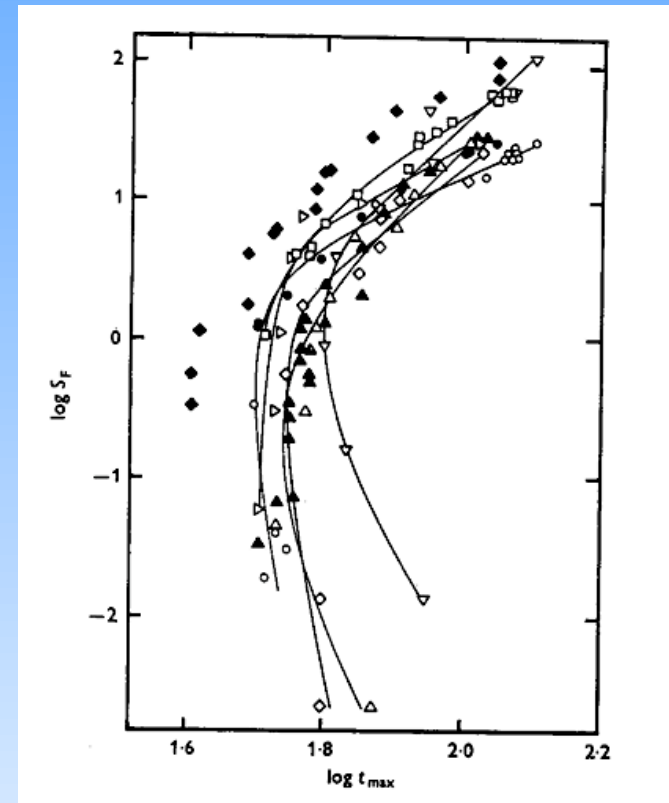


(Baylor and Hodgkin, 1974)

# $I_{\text{pred}}$ extraction and maximization

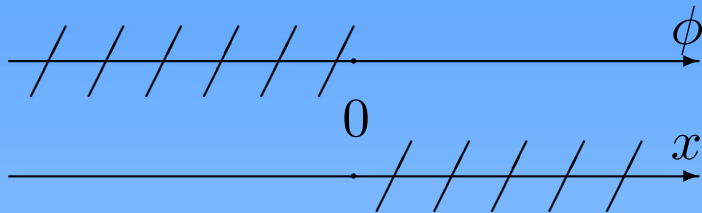


$I([x_{\text{past}}], [\phi_{\text{future}}])$  – too difficult



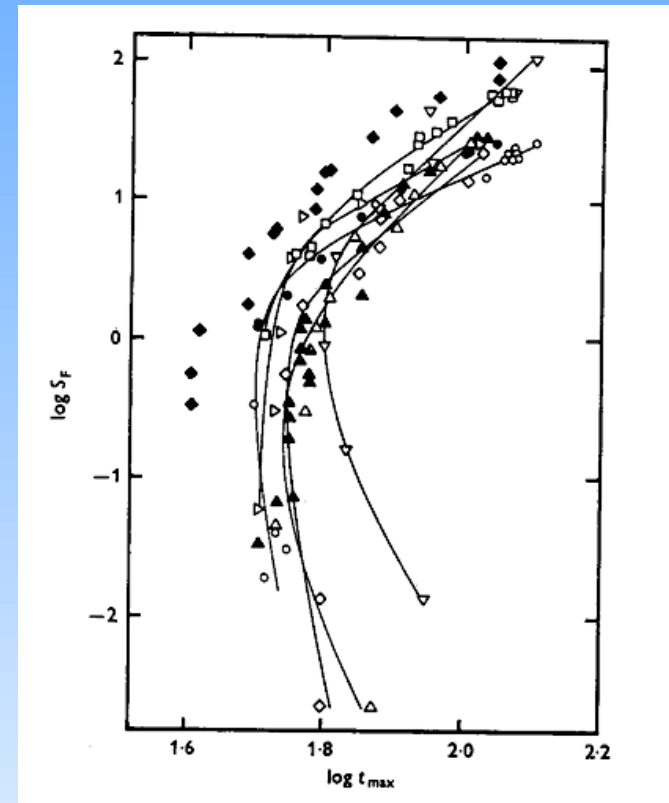
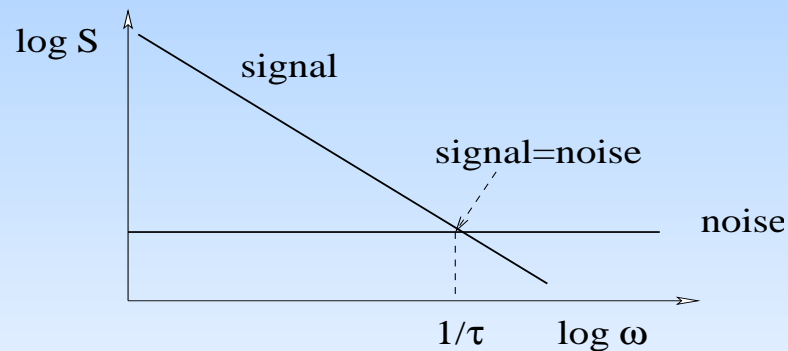
(Baylor and Hodgkin, 1974)

# $I_{\text{pred}}$ extraction and maximization



$I([x_{\text{past}}], [\phi_{\text{future}}])$  – too difficult

$$I(x_0, \phi_0) = \log \frac{\langle \phi^2 \rangle}{\langle \phi^2 \rangle - \frac{\langle \phi_f^2 \rangle^2}{\langle x^2 \rangle}}$$



(Baylor and Hodgkin, 1974)

## Specific examples: problem setup

$Q(\vec{x}|\alpha)$  p. d. f. for  $\vec{x}$  parameterized by unknown parameters  $\alpha$

$\dim \alpha = K$  dimensionality of  $\alpha$ , may be infinite

$\mathcal{P}(\alpha)$  prior distribution of parameters

$\vec{x}_1 \cdots \vec{x}_N$  random samples from the distribution

## Specific examples: problem setup

$Q(\vec{x}|\alpha)$  p. d. f. for  $\vec{x}$  parameterized by unknown parameters  $\alpha$

$\dim \alpha = K$  dimensionality of  $\alpha$ , may be infinite

$\mathcal{P}(\alpha)$  prior distribution of parameters

$\vec{x}_1 \cdots \vec{x}_N$  random samples from the distribution

$$P(\vec{x}_1, \vec{x}_2, \cdots, \vec{x}_N | \alpha) = \prod_{i=1}^N Q(\vec{x}_i | \alpha)$$

$$P(\vec{x}_1, \vec{x}_2, \cdots, \vec{x}_N) = \int d^K \alpha \mathcal{P}(\alpha) \prod_{i=1}^N Q(\vec{x}_i | \alpha)$$

$$S(\vec{x}_1, \vec{x}_2, \cdots, \vec{x}_N) \equiv S(N)$$

$$= - \int d\vec{x}_1 \cdots d\vec{x}_N P(\{\vec{x}_i\}) \log_2 P(\{\vec{x}_i\})$$



## Separating the terms

$$\mathcal{S}_0 = \int d^K \alpha \mathcal{P}(\alpha) \left[ - \int d\vec{x} Q(\vec{x}|\alpha) \log_2 Q(\vec{x}|\alpha) \right]$$

$$\mathcal{S}_1(N) = - \int d^K \bar{\alpha} d^N \vec{x}_i \mathcal{P}(\bar{\alpha}) \prod Q(\vec{x}_i|\bar{\alpha}) \log_2 \int d^K \alpha \mathcal{P}(\alpha) e^{-N\mathcal{E}_N}$$

## Separating the terms

$$\mathcal{S}_0 = \int d^K \alpha \mathcal{P}(\alpha) \left[ - \int d\vec{x} Q(\vec{x}|\alpha) \log_2 Q(\vec{x}|\alpha) \right]$$

$$\mathcal{S}_1(N) = - \int d^K \bar{\alpha} d^N \vec{x}_i \mathcal{P}(\bar{\alpha}) \prod Q(\vec{x}_i|\bar{\alpha}) \log_2 \int d^K \alpha \mathcal{P}(\alpha) e^{-N\mathcal{E}_N}$$

$$\mathcal{E}_N \equiv \frac{1}{N} \sum_i \log \left[ \frac{Q(\vec{x}_i|\bar{\alpha})}{Q(\vec{x}_i|\alpha)} \right] \xrightarrow{\text{anneal}} \int d\vec{x} Q(\vec{x}|\bar{\alpha}) \log \frac{Q(\vec{x}|\bar{\alpha})}{Q(\vec{x}|\alpha)}$$

Annealed approximation (almost) always works.

## Density of states

$$Z(\bar{\alpha}; N) = \int d\epsilon \rho(\epsilon; \bar{\alpha}) \exp[-N\epsilon]$$

$$\rho(\epsilon; \bar{\alpha}) = \int d^K \alpha \mathcal{P}(\alpha) \delta[\epsilon - D_{\text{KL}}(\bar{\alpha} || \alpha)]$$

## Density of states

$$Z(\bar{\alpha}; N) = \int d\epsilon \rho(\epsilon; \bar{\alpha}) \exp[-N\epsilon]$$

$$\rho(\epsilon; \bar{\alpha}) = \int d^K \alpha \mathcal{P}(\alpha) \delta[\epsilon - D_{\text{KL}}(\bar{\alpha} || \alpha)]$$

$$\int d\epsilon \rho(\epsilon; \bar{\alpha}) = \int d^K \alpha \mathcal{P}(\alpha) = 1 \quad \text{annealing works!}$$

## Density of states

$$Z(\bar{\alpha}; N) = \int d\epsilon \rho(\epsilon; \bar{\alpha}) \exp[-N\epsilon]$$

$$\rho(\epsilon; \bar{\alpha}) = \int d^K \alpha \mathcal{P}(\alpha) \delta[\epsilon - D_{\text{KL}}(\bar{\alpha} || \alpha)]$$

$$\int d\epsilon \rho(\epsilon; \bar{\alpha}) = \int d^K \alpha \mathcal{P}(\alpha) = 1 \quad \text{annealing works!}$$

Learning is annealing at decreasing temperature.

Nonzero  $\rho \implies$  consistency in learning.

## Density at $\epsilon \rightarrow 0$ , $I_{\text{pred}}$ , and learning

Occam factor, generalization error, prediction error, fluctuation determinant:

$$\mathcal{D}(\bar{\alpha}; N) \approx -\log \int d\epsilon \rho(\epsilon; \bar{\alpha}) e^{-N\epsilon}$$

## Density at $\epsilon \rightarrow 0$ , $I_{\text{pred}}$ , and learning

Occam factor, generalization error, prediction error, fluctuation determinant:

$$\mathcal{D}(\bar{\alpha}; N) \approx -\log \int d\epsilon \rho(\epsilon; \bar{\alpha}) e^{-N\epsilon}$$

Predictive information:

$$I_{\text{pred}}(N) \approx \int d^K \bar{\alpha} \mathcal{P}(\bar{\alpha}) \mathcal{D}(\bar{\alpha}, N)$$

## Density at $\epsilon \rightarrow 0$ , $I_{\text{pred}}$ , and learning

Occam factor, generalization error, prediction error, fluctuation determinant:

$$\mathcal{D}(\bar{\alpha}; N) \approx -\log \int d\epsilon \rho(\epsilon; \bar{\alpha}) e^{-N\epsilon}$$

Predictive information:

$$I_{\text{pred}}(N) \approx \int d^K \bar{\alpha} \mathcal{P}(\bar{\alpha}) \mathcal{D}(\bar{\alpha}, N)$$

Universal learning curves:

$$\Lambda(\bar{\alpha}; N) \equiv D_{\text{KL}}(\bar{\alpha} || \alpha_{\text{est}}) \approx \frac{d\mathcal{D}(\bar{\alpha}; N)}{dN}$$

$$\Lambda(N) \equiv \int d\bar{\alpha} \mathcal{P}(\bar{\alpha}) \Lambda(\bar{\alpha}; N) \approx \frac{dI_{\text{pred}}}{dN}$$



# Finite number of states and finite $I_{\text{pred}}$

$$\rho(\epsilon; a_i) = \sum_{j=1}^M \mathcal{P}_j \delta(d_{ij} - \epsilon)$$

# Finite number of states and finite $I_{\text{pred}}$

$$\rho(\epsilon; a_i) = \sum_{j=1}^M \mathcal{P}_j \delta(d_{ij} - \epsilon)$$

$$\mathcal{D}(a_i; N) = c_1 - c_2 \exp[-Nc_3]$$

$$\Lambda(a_i; N) \approx c_2 c_3 \exp[-Nc_3]$$

$I_{\text{pred}}$  saturates as  $N \rightarrow \infty$

# Power-law density function

$$\rho(\epsilon \rightarrow 0; \bar{\alpha}) \approx A(\bar{\alpha}) \epsilon^{(d-2)/2}$$

## Power-law density function

$$\rho(\epsilon \rightarrow 0; \bar{\alpha}) \approx A(\bar{\alpha}) \epsilon^{(d-2)/2}$$

**Example:** *sound* finite parameter models,  $\dim \alpha = d$ .

$$\rho(\epsilon; \bar{\alpha}) \xrightarrow{\epsilon \rightarrow 0} \mathcal{P}(\bar{\alpha}) \frac{2\pi^{d/2}}{\Gamma(d/2)} (\det \mathcal{F})^{-1/2} \epsilon^{(d-2)/2}$$

$$I_{\text{pred}} \approx S_1^{(\text{a})} \approx \frac{d}{2} \log_2 N$$

## Power-law density function

$$\rho(\epsilon \rightarrow 0; \bar{\alpha}) \approx A(\bar{\alpha}) \epsilon^{(d-2)/2}$$

**Example:** *sound* finite parameter models,  $\dim \alpha = d$ .

$$\rho(\epsilon; \bar{\alpha}) \xrightarrow{\epsilon \rightarrow 0} \mathcal{P}(\bar{\alpha}) \frac{2\pi^{d/2}}{\Gamma(d/2)} (\det \mathcal{F})^{-1/2} \epsilon^{(d-2)/2}$$

$$I_{\text{pred}} \approx S_1^{(a)} \approx \frac{d}{2} \log_2 N$$

Speed of approach to this asymptotics is rarely investigated.

## Another example

Learning  $Q(\vec{x}_1 \cdots \vec{x}_N | \alpha)$ , a finite parameter Markov process with long range intrinsic correlations such that

$$\begin{aligned} S[\{\vec{x}_i\} | \alpha] &\equiv - \int d^N \vec{x} Q(\{\vec{x}_i\} | \alpha) \log_2 Q(\{\vec{x}_i\} | \alpha) \\ &\rightarrow N \mathcal{S}_0 + \mathcal{S}_0^*; \quad \mathcal{S}_0^* = \frac{K'}{2} \log_2 N \end{aligned}$$

## Another example

Learning  $Q(\vec{x}_1 \cdots \vec{x}_N | \alpha)$ , a finite parameter Markov process with long range intrinsic correlations such that

$$S[\{\vec{x}_i\} | \alpha] \equiv - \int d^N \vec{x} Q(\{\vec{x}_i\} | \alpha) \log_2 Q(\{\vec{x}_i\} | \alpha)$$

$$\rightarrow N\mathcal{S}_0 + \mathcal{S}_0^*; \quad \mathcal{S}_0^* = \frac{K'}{2} \log_2 N$$

$$S_1^{(a)}(N) \approx \frac{K + K'}{2} \log_2 N$$

## Another example

Learning  $Q(\vec{x}_1 \cdots \vec{x}_N | \alpha)$ , a finite parameter Markov process with long range intrinsic correlations such that

$$S[\{\vec{x}_i\} | \alpha] \equiv - \int d^N \vec{x} Q(\{\vec{x}_i\} | \alpha) \log_2 Q(\{\vec{x}_i\} | \alpha)$$

$$\rightarrow N\mathcal{S}_0 + \mathcal{S}_0^*; \quad \mathcal{S}_0^* = \frac{K'}{2} \log_2 N$$

$$S_1^{(a)}(N) \approx \frac{K + K'}{2} \log_2 N$$

Do not distinguish predictability from unknown parameters and from intrinsic correlations.

In physics similar to: order parameters  $\Longleftrightarrow$  interactions.



## Essential singularity in the density

$$\rho(\epsilon \rightarrow 0; \bar{\alpha}) \approx A(\bar{\alpha}) \exp \left[ -\frac{B(\bar{\alpha})}{\epsilon^\mu} \right], \quad \mu > 0$$

$$S_1^{(a)}(N) \propto N^{\mu/(\mu+1)}$$

## Essential singularity in the density

$$\rho(\epsilon \rightarrow 0; \bar{\alpha}) \approx A(\bar{\alpha}) \exp \left[ -\frac{B(\bar{\alpha})}{\epsilon^\mu} \right], \quad \mu > 0$$

$$S_1^{(a)}(N) \propto N^{\mu/(\mu+1)}$$

- finite parameter model with increasing number of parameters  $K \sim N^{\mu/(\mu+1)}$ ;  $S_1(N) \sim N^{\mu/(\mu+1)}$ , not  $S_1(N) \sim \frac{N^{\mu/(\mu+1)}}{2} \log N$
- as  $\mu \rightarrow \infty$  complexity grows and then vanishes to the leading order when  $S_1^{(a)}$  becomes extensive

## Example of the power-law $I_{\text{pred}}$

Learning a smooth nonparameteric density  $Q(x) = 1/l_0 e^{-\phi(x)}$ ,  
 $x \in [0, L]$  (Bialek, Callan, and Strong 1996), **Complete model.**

$$\mathcal{P}[\phi(x)] = \frac{1}{\mathcal{Z}} \exp \left[ -\frac{l}{2} \int dx \left( \frac{\partial \phi}{\partial x} \right)^2 \right] \delta \left[ \frac{1}{l_0} \int dx e^{-\phi(x)} - 1 \right]$$

## Example of the power-law $I_{\text{pred}}$

Learning a smooth nonparameteric density  $Q(x) = 1/l_0 e^{-\phi(x)}$ ,  
 $x \in [0, L]$  (Bialek, Callan, and Strong 1996), **Complete model.**

$$\mathcal{P}[\phi(x)] = \frac{1}{\mathcal{Z}} \exp \left[ -\frac{l}{2} \int dx \left( \frac{\partial \phi}{\partial x} \right)^2 \right] \delta \left[ \frac{1}{l_0} \int dx e^{-\phi(x)} - 1 \right]$$

$$\rho(D \rightarrow 0; \bar{\phi}) = A[\bar{\phi}(x)] \epsilon^{-3/2} \exp \left( -\frac{B[\bar{\phi}(x)]}{\epsilon} \right)$$

$$S_1^{(a)}(N) \propto \sqrt{N} \left( \frac{L}{l} \right)^{1/2}$$

## Example of the power-law $I_{\text{pred}}$

Learning a smooth nonparameteric density  $Q(x) = 1/l_0 e^{-\phi(x)}$ ,  
 $x \in [0, L]$  (Bialek, Callan, and Strong 1996), **Complete model**.

$$\mathcal{P}[\phi(x)] = \frac{1}{\mathcal{Z}} \exp \left[ -\frac{l}{2} \int dx \left( \frac{\partial \phi}{\partial x} \right)^2 \right] \delta \left[ \frac{1}{l_0} \int dx e^{-\phi(x)} - 1 \right]$$

$$\rho(D \rightarrow 0; \bar{\phi}) = A[\bar{\phi}(x)] \epsilon^{-3/2} \exp \left( -\frac{B[\bar{\phi}(x)]}{\epsilon} \right)$$

$$S_1^{(a)}(N) \propto \sqrt{N} \left( \frac{L}{l} \right)^{1/2}$$

- increasing number of “effective parameters” (bins) of adaptive size  $\sim \sqrt{l/NQ(x)}$

## Example of the power-law $I_{\text{pred}}$

Learning a smooth nonparameteric density  $Q(x) = 1/l_0 e^{-\phi(x)}$ ,  
 $x \in [0, L]$  (Bialek, Callan, and Strong 1996), **Complete model**.

$$\mathcal{P}[\phi(x)] = \frac{1}{\mathcal{Z}} \exp \left[ -\frac{l}{2} \int dx \left( \frac{\partial \phi}{\partial x} \right)^2 \right] \delta \left[ \frac{1}{l_0} \int dx e^{-\phi(x)} - 1 \right]$$

$$\rho(D \rightarrow 0; \bar{\phi}) = A[\bar{\phi}(x)] \epsilon^{-3/2} \exp \left( -\frac{B[\bar{\phi}(x)]}{\epsilon} \right)$$

$$S_1^{(a)}(N) \propto \sqrt{N} \left( \frac{L}{l} \right)^{1/2}$$

- increasing number of “effective parameters” (bins) of adaptive size  $\sim \sqrt{l/NQ(x)}$
- heuristic arguments for the dimensionality  $\zeta$  and the smoothness exponent  $\eta$  give  $S_1(N) \sim N^{\zeta/2\eta}$  — demonstrates a crossover from complexity to randomness

## Nonuniform $\mathcal{D}$

Nested finite parameter models,  $r = 1 \dots \infty$ ,  $K = K(r)$ ,  $\mathcal{P}(r)$ :

$$\mathcal{P}(\alpha_\mu|r) = \begin{cases} p(\alpha_\mu), & \mu \leq K(r) \\ \delta(\alpha_\mu), & \mu > K(r) \end{cases}$$

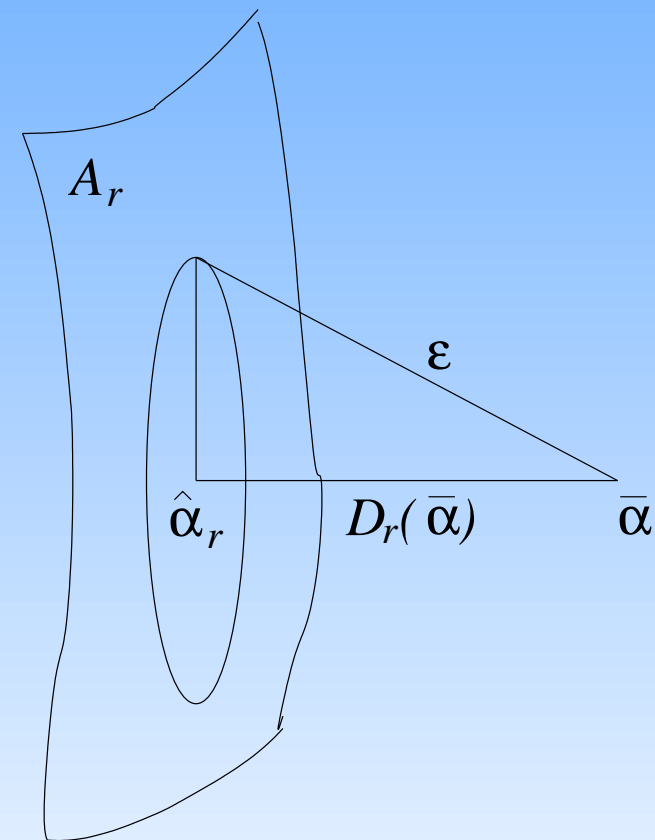
$$\mathcal{P}(\boldsymbol{\alpha}|r) = \prod_{\mu=1}^R \mathcal{P}(\alpha_\mu|r)$$

# Nonuniform $\mathcal{D}$

Nested finite parameter models,  $r = 1 \dots \infty$ ,  $K = K(r)$ ,  $\mathcal{P}(r)$ :

$$\mathcal{P}(\alpha_\mu|r) = \begin{cases} p(\alpha_\mu), & \mu \leq K(r) \\ \delta(\alpha_\mu), & \mu > K(r) \end{cases}$$

$$\mathcal{P}(\boldsymbol{\alpha}|r) = \prod_{\mu=1}^R \mathcal{P}(\alpha_\mu|r)$$





# Nonuniform $\mathcal{D}$

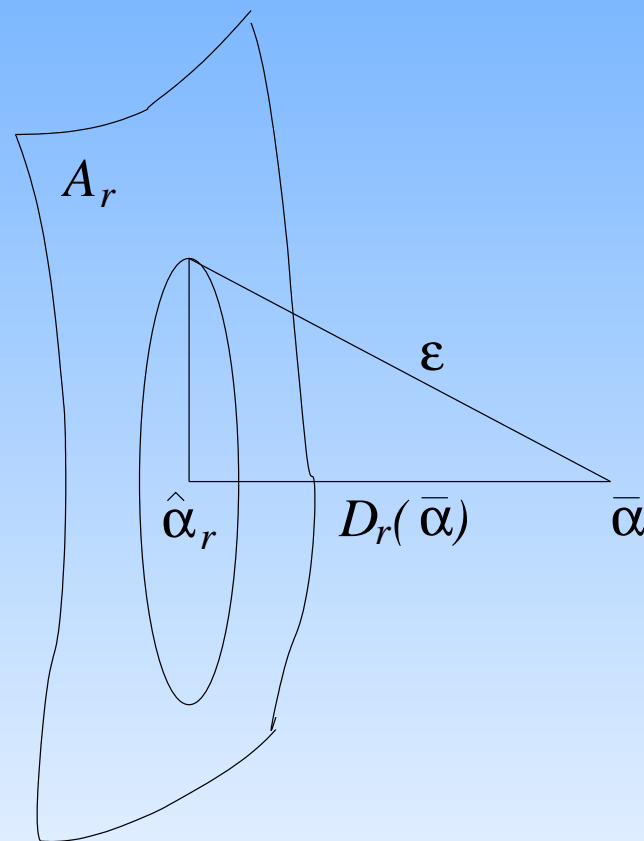
Nested finite parameter models,  $r = 1 \dots \infty$ ,  $K = K(r)$ ,  $\mathcal{P}(r)$ :

$$\mathcal{P}(\alpha_\mu|r) = \begin{cases} p(\alpha_\mu), & \mu \leq K(r) \\ \delta(\alpha_\mu), & \mu > K(r) \end{cases}$$

$$\mathcal{P}(\boldsymbol{\alpha}|r) = \prod_{\mu=1}^R \mathcal{P}(\alpha_\mu|r)$$

$$\rho(\epsilon; \bar{\boldsymbol{\alpha}}) = \sum_{r: D_r(\bar{\boldsymbol{\alpha}}) \leq \epsilon} \mathcal{P}(r) \mathcal{P}(\hat{\boldsymbol{\alpha}}_r|r)$$

$$\frac{2\pi^{K(r)/2}}{\Gamma[K(r)/2]} \frac{[\epsilon^2 - D_r^2(\bar{\boldsymbol{\alpha}})]^{[K(r)-2]/4}}{\sqrt{\det \mathcal{F}_{K(r)}}}$$



# Nonuniform $\mathcal{D}$

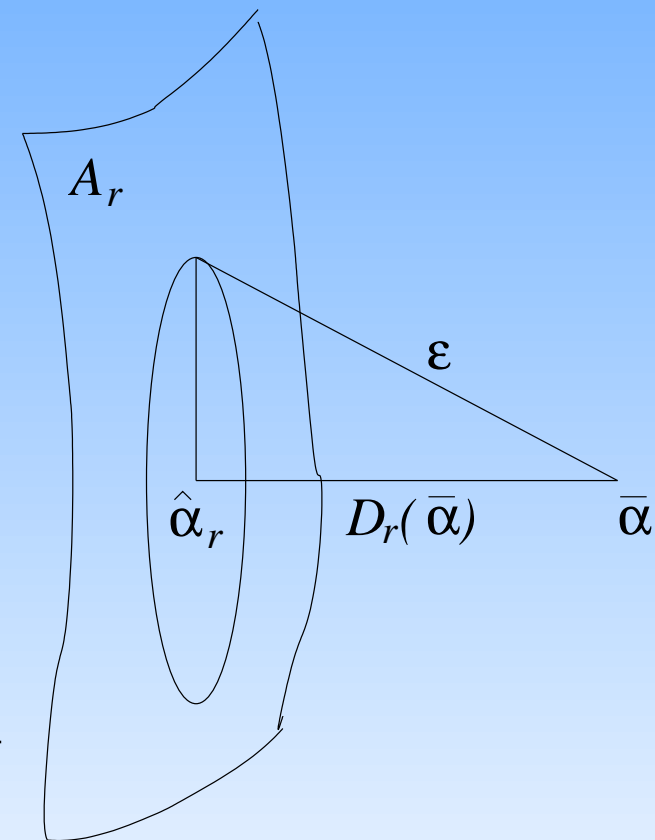
Nested finite parameter models,  $r = 1 \dots \infty$ ,  $K = K(r)$ ,  $\mathcal{P}(r)$ :

$$\mathcal{P}(\alpha_\mu|r) = \begin{cases} p(\alpha_\mu), & \mu \leq K(r) \\ \delta(\alpha_\mu), & \mu > K(r) \end{cases}$$

$$\mathcal{P}(\boldsymbol{\alpha}|r) = \prod_{\mu=1}^R \mathcal{P}(\alpha_\mu|r)$$

$$\rho(\epsilon; \bar{\boldsymbol{\alpha}}) = \sum_{r: D_r(\bar{\boldsymbol{\alpha}}) \leq \epsilon} \mathcal{P}(r) \mathcal{P}(\hat{\boldsymbol{\alpha}}_r|r)$$

$$\frac{2\pi^{K(r)/2}}{\Gamma[K(r)/2]} \frac{[\epsilon^2 - D_r^2(\bar{\boldsymbol{\alpha}})]^{[K(r)-2]/4}}{\sqrt{\det \mathcal{F}_{K(r)}}}$$



Another complete model!

## Density near atypical solutions

$$\rho(\epsilon; \bar{\alpha}) \sim \epsilon^{\epsilon^{-1/(2\eta-1)} \ell^{-1}}$$

$$\mathcal{D} \propto N^{1/2\eta} \left( \frac{\log N}{\ell} \right)^{1-1/2\eta}$$

## Density near atypical solutions

$$\rho(\epsilon; \bar{\alpha}) \sim \epsilon^{\epsilon^{-1/(2\eta-1)} \ell^{-1}}$$

$$\mathcal{D} \propto N^{1/2\eta} \left( \frac{\log N}{\ell} \right)^{1-1/2\eta}$$

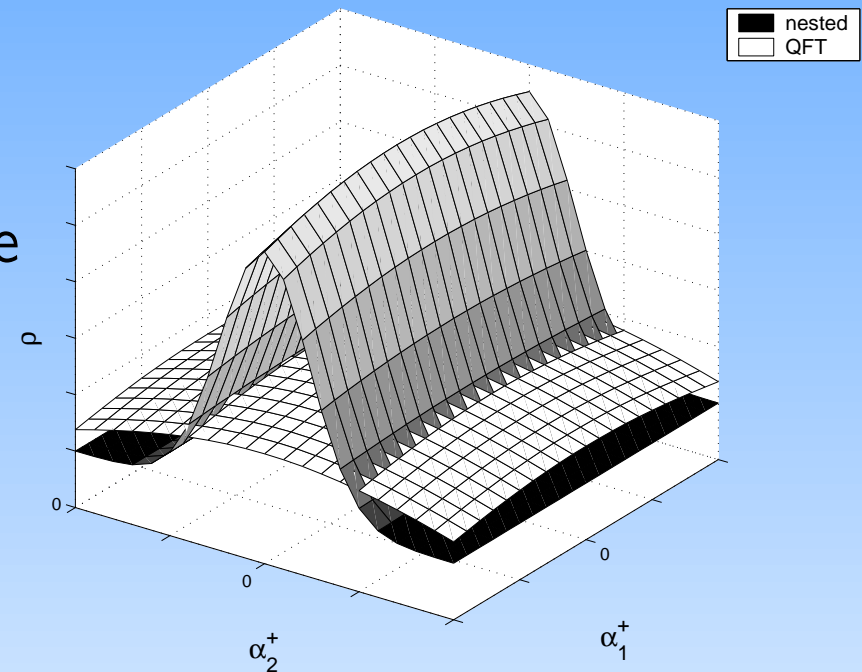
- nested model is at most log worse than the QFT
- QFT may be a power law worse

## Density near atypical solutions

$$\rho(\epsilon; \bar{\alpha}) \sim \epsilon^{\epsilon^{-1/(2\eta-1)} \ell^{-1}}$$

$$\mathcal{D} \propto N^{1/2\eta} \left( \frac{\log N}{\ell} \right)^{1-1/2\eta}$$

- nested model is at most log worse than the QFT
- QFT may be a power law worse

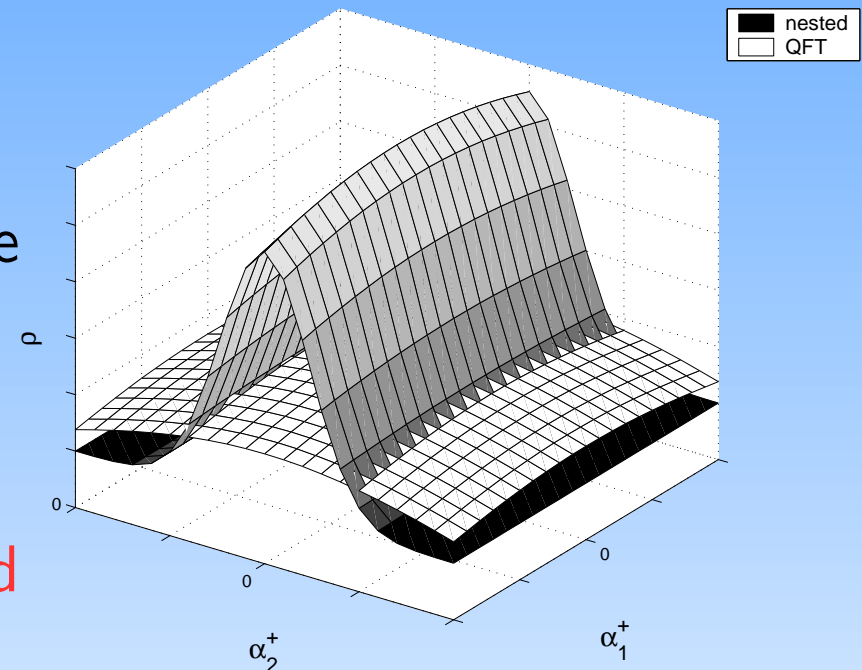


## Density near atypical solutions

$$\rho(\epsilon; \bar{\alpha}) \sim \epsilon^{\epsilon^{-1/(2\eta-1)} \ell^{-1}}$$

$$\mathcal{D} \propto N^{1/2\eta} \left( \frac{\log N}{\ell} \right)^{1-1/2\eta}$$

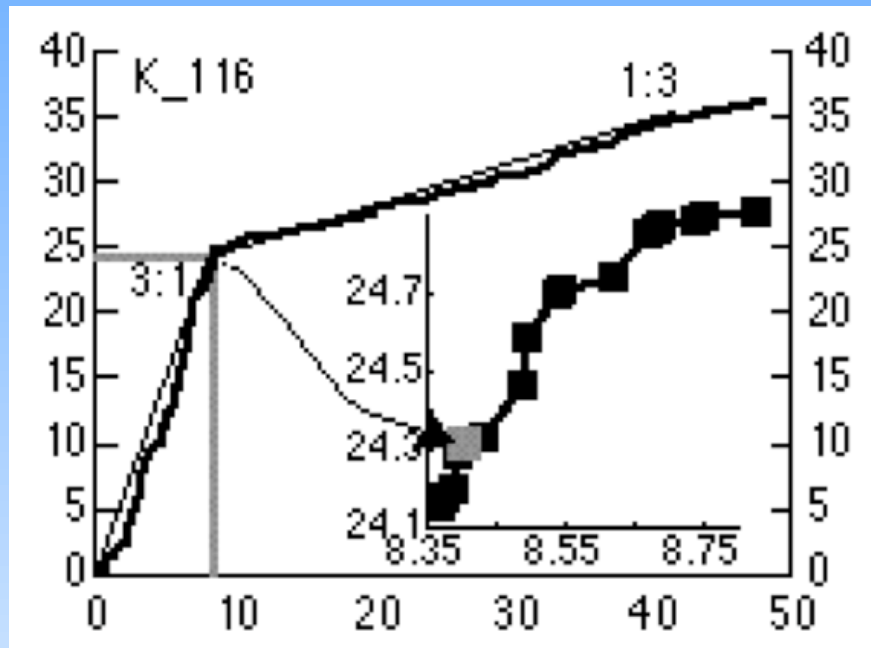
- nested model is at most log worse than the QFT
- QFT may be a power law worse
- for natural (structured) data nested case is better
- alignment may be imperfect for finite precision  $\epsilon$



# Which model is being used?

- for QFT or nested asymptotics  
kicks in fast
- asymptotic decay rate should  
signify the model

## Which model is being used?

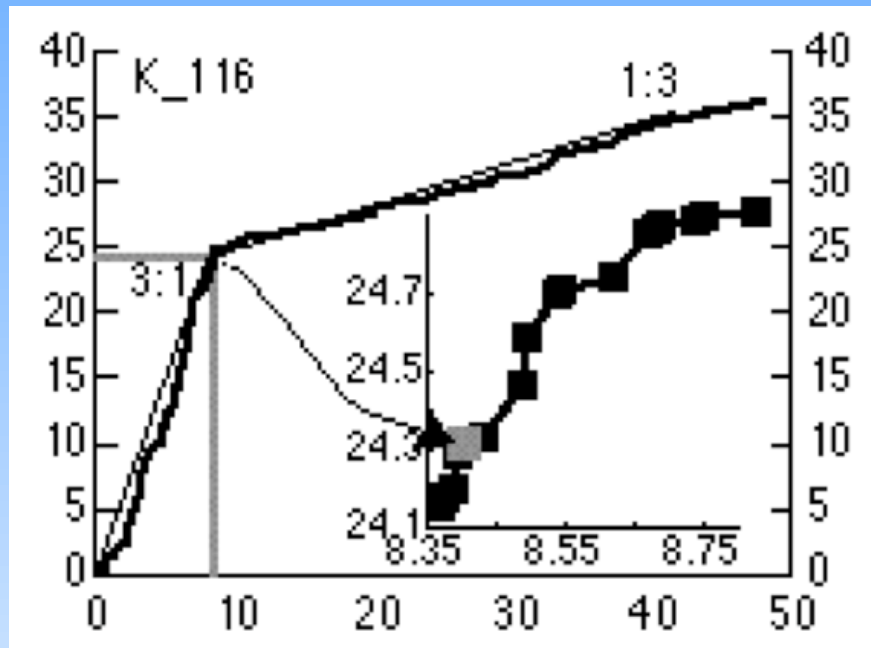


- for QFT or nested asymptotics kicks in fast
- asymptotic decay rate should signify the model
- decay rate too fast to observe
- noisy learning

(Gallistel et al., 2001)



# Which model is being used?

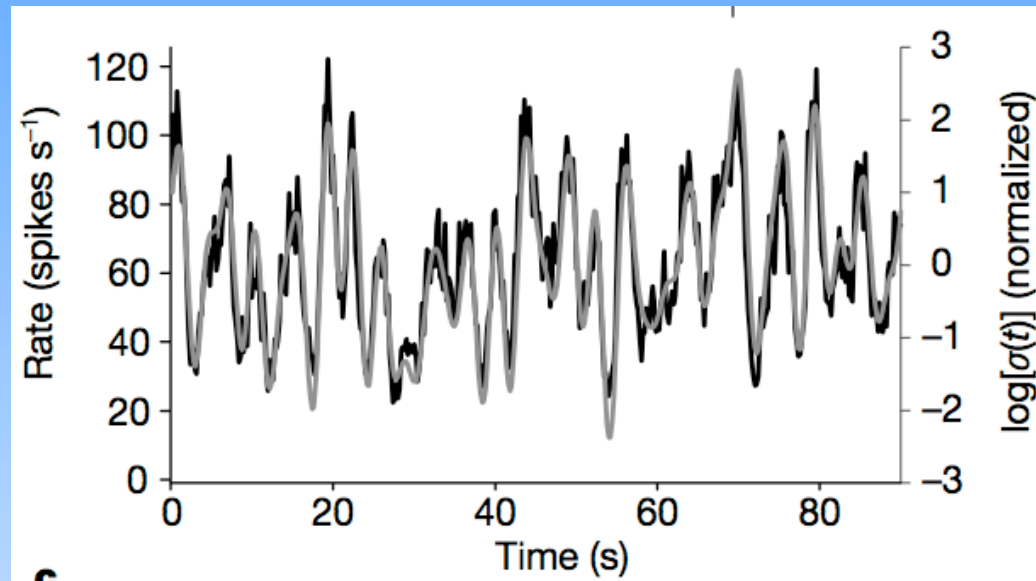


- for QFT or nested asymptotics kicks in fast
- asymptotic decay rate should signify the model
- decay rate too fast to observe
- noisy learning

- maybe FDT?  $\frac{\partial \Lambda}{\partial N} = -\zeta_N \Lambda^\nu$

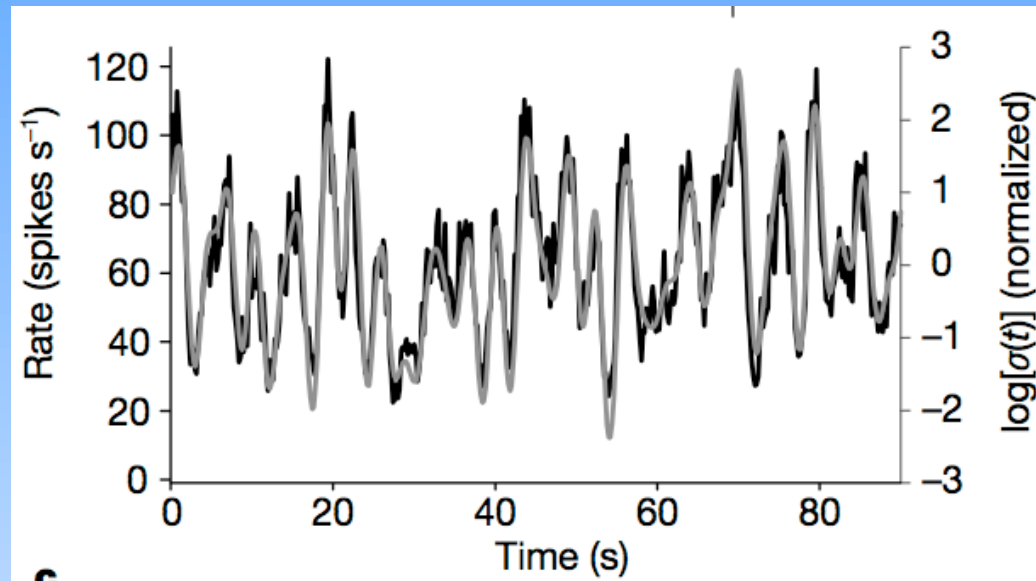
(Gallistel et al., 2001)

# Fluctuations (drifting target) and dissipation (learning curve)



(Fairhall et al., 2001)

# Fluctuations (drifting target) and dissipation (learning curve)



(Fairhall et al., 2001)

$$\Delta_{\text{rms}} = \left\{ \nu^{1/\nu} \frac{\Gamma\left(\frac{3}{2\nu}\right)}{\Gamma\left(\frac{1}{2\nu}\right)} \right\}^{1/2} \left( \frac{\Omega}{\zeta} \right)^{1/(2\nu)}$$

# The hidden extras. . .

# Relations to other definitions of complexity ...

... are mostly straightforward.

## Relations to other definitions of complexity ...

... are mostly straightforward.

For Kolmogorov complexity:

## Relations to other definitions of complexity ...

... are mostly straightforward.

For Kolmogorov complexity:

- partition all strings into equivalence classes

## Relations to other definitions of complexity ...

... are mostly straightforward.

For Kolmogorov complexity:

- partition all strings into equivalence classes
- define Kolmogorov complexity  $C_K(s)$  of a sequence  $s$  with respect to the partition as a length of the shortest program that can generate a sequence from the class  $s$  belongs to



## Relations to other definitions of complexity ...

... are mostly straightforward.

For Kolmogorov complexity:

- partition all strings into equivalence classes
- define Kolmogorov complexity  $C_K(s)$  of a sequence  $s$  with respect to the partition as a length of the shortest program that can generate a sequence from the class  $s$  belongs to
- equivalence = indistinguishable conditional distributions of futures

## Relations to other definitions of complexity ...

... are mostly straightforward.

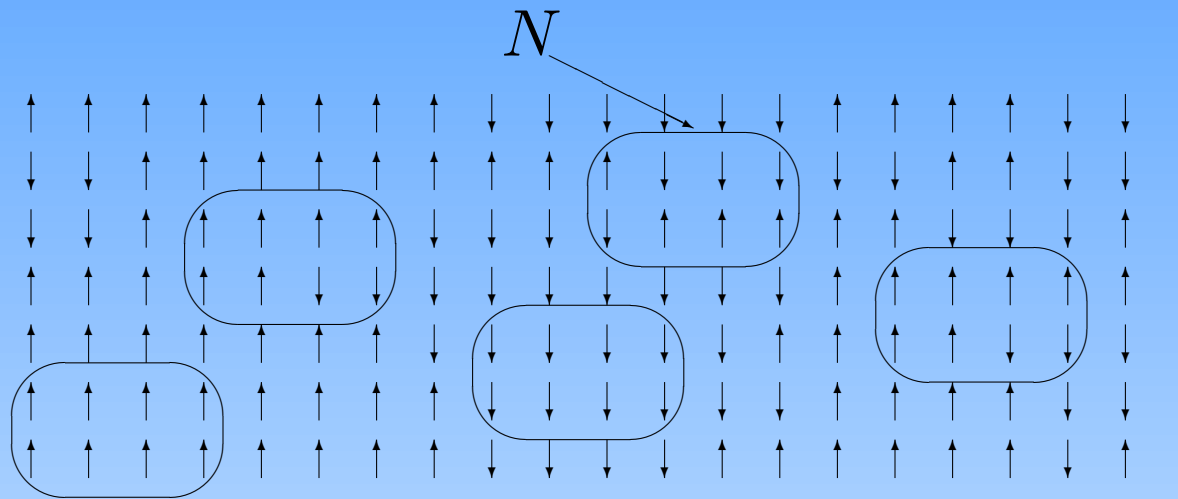
For Kolmogorov complexity:

- partition all strings into equivalence classes
- define Kolmogorov complexity  $C_K(s)$  of a sequence  $s$  with respect to the partition as a length of the shortest program that can generate a sequence from the class  $s$  belongs to
- equivalence = indistinguishable conditional distributions of futures

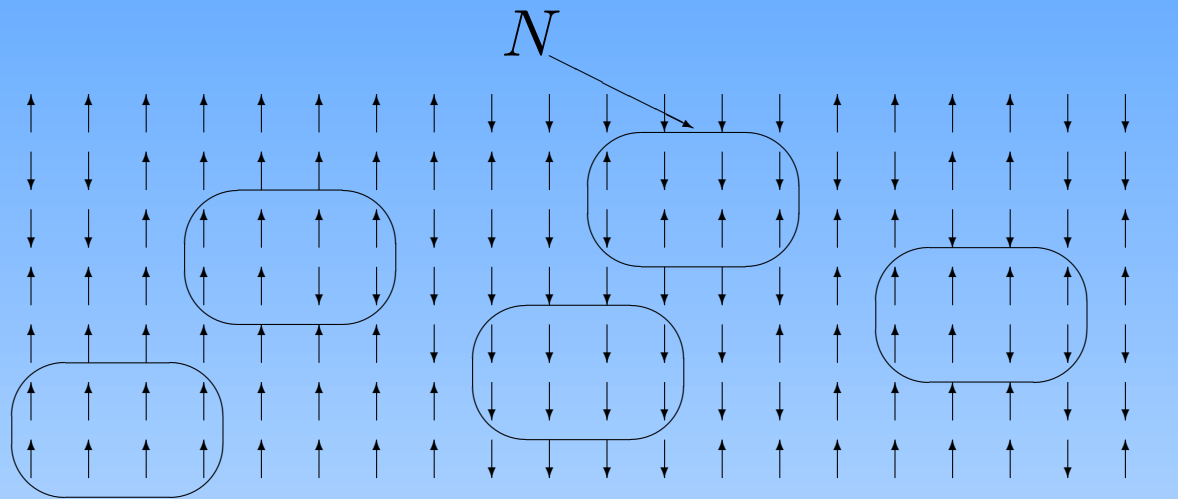
If sufficient statistics exist, then  $C_K \approx I_{\text{pred}}$ . Otherwise  $C_K > I_{\text{pred}}$ .

$C_K$  is unique up to a constant.

# RG, not finite size scaling!



# RG, not finite size scaling!



$$S(N) = S(\text{block}) + S(\text{spin}|\text{block})$$

Scaling fields carry information across.

Is  $I_{\text{pred}} = f(\text{scaling exponents}) \log N$ ?

## What's next?

**extraction** separating predictive information from non-predictive using the *Information Bottleneck* technique

**physics** of phase transitions, connection to subextensive statistical mechanics

**learning** unification of approaches: Bayesian, SRM, MDL, Cucker-Smale. . .

**bioinformatics** what is predictive information of natural symbolic sequences (DNA, languages, spike trains)? animal behavior? can we use changes in predictability for data partitioning? for model building?

**dynamical systems theory** what is predictive information and complexity of various systems?