

# Neural dynamics of vision

## A computational perspective

CLAYTON SEITZ<sup>1</sup>

February 27, 2021



Dedicated to Calvin and Hobbes.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Natural Image Statistics</b>	<b>3</b>
2.1	Image Statistics . . . . .	3
2.2	Section heading . . . . .	3
<b>3</b>	<b>Information and Coding Theory</b>	<b>5</b>
3.1	Introduction . . . . .	5
3.2	Entropy . . . . .	5
3.2.1	Joint and Conditional Entropy . . . . .	6
3.3	KL-Divergence and Mutual Information . . . . .	7
3.3.1	The Data-Processing Inequality . . . . .	8
3.4	Source Coding . . . . .	8
3.4.1	Kraft's Inequality . . . . .	9
3.4.2	Source Coding Theorem . . . . .	9
3.5	Error Correcting Codes (Channel Coding) . . . . .	9
<b>4</b>	<b>The Neural Code</b>	<b>11</b>
4.1	Introduction . . . . .	11
4.2	Neuroelectronics . . . . .	11
4.2.1	RC Model . . . . .	12
4.2.2	Equilibrium and Reversal Potentials . . . . .	13
4.2.3	Membrane Current . . . . .	14



# 1

## Introduction

The purpose of the following chapters is to highlight some of the statistical properties of visual inputs and neural responses. In particular, we will discuss the statistics and information theory of natural images and how constraints on the space of inputs can in turn constrain the response statistics. Then, we will discuss some of the ways the human visual system might compare inputs to this learned model - a phenomenon we will call *visual alignment*.





## 2

# Natural Image Statistics

The study of natural images has recently taken hold in the neuroscience community and there are several reasons for this, one of which is that the human visual system and visual systems in other organisms have developed in an environment sampled by these images. It is then reasonable to hypothesize that these visual systems have developed within and optimized for these kinds of inputs. To be more precise, the architecture of the visual system could be wired to represent the statistical distribution of natural scenes. On the other hand, natural images are a very sparse subset of the entire space of possible images - natural images are speckled throughout image space according to some kind of natural image distribution. The volume of the space of possible images consumed by natural images is unimaginably small but we still might ask, is it possible to write down features of the natural image distribution over the space of images? In other words, we would like to determine the shape of the distribution, be it isolated pockets in the space, a hyperplane, etc. Then, even if we could write down such features, how would a network of neurons encode it?

## 2.1 Image Statistics

The treatment used throughout this chapter will draw heavily from information theory, coding theory, and statistics. It is by no means a comprehensive introduction to any of these topics, although many of the fundamental concepts will be introduced. By the end of this chapter, we will have a better idea of how to formulate a generative model, which produces a natural image from a set of latent variables.

## 2.2 Section heading

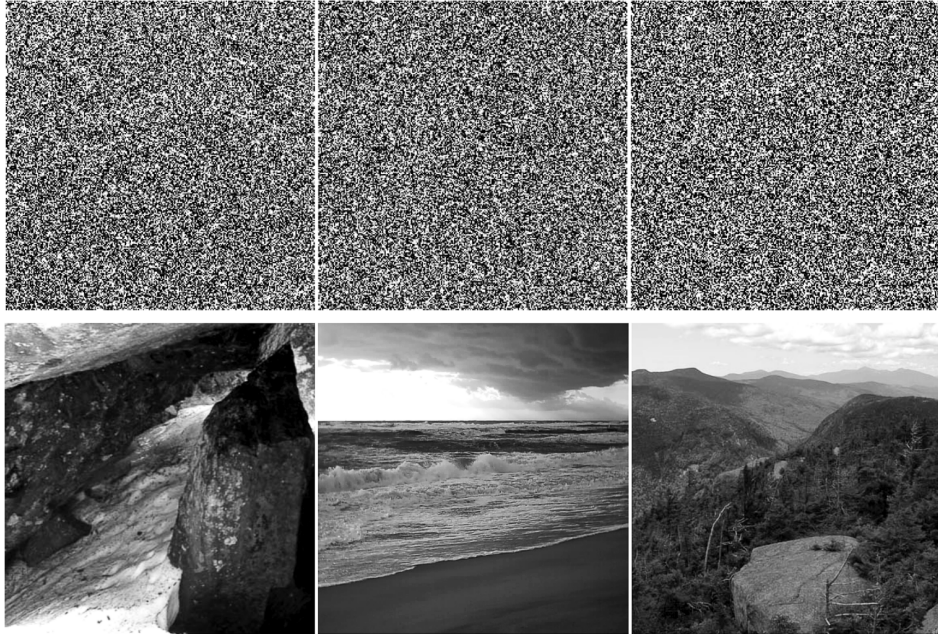


Figure 2.1: Comparison of images generated from white noise with natural images

# 3

## Information and Coding Theory

*“We may have knowledge of the past but cannot control it; we may control the future but have no knowledge of it”*

– Claude Shannon

### 3.1 Introduction

Information theory is a framework first introduced by Claude Shannon’s seminal paper *A mathematical theory of communication* published in 1948. At its core, information theory makes the intuitive concept of *information* mathematically rigorous and forms the foundation of many modern communication systems. Neural circuits in the visual system are an especially interesting example of such a communication system. Therefore, in this section, the information theoretic concepts necessary for studying neural circuits are introduced.

### 3.2 Entropy

The concept of entropy is not exclusive to information theory; rather, it is used widely in disciplines such as physics and mathematical statistics. In fact, entropy was originally defined in statistical physics when Ludwig Boltzmann gave a statistical description of a thermodynamic system of particles. Since this is arguably the more intuitive path as opposed to an entirely mathematical description, I will follow a similar line of reasoning in the following paragraphs.

In every application, the entropy  $\mathbf{H}$  is a measure of uncertainty or how much information is contained in a random variable  $x$ . In information theory, the entropy is a property of a probability distribution of a random variable

$P(x)$  where  $x$  can take on continuous or discrete values. For the discrete case, we can express the entropy in bits

$$\mathbf{H} = \sum_{x \in S} P(x) \log \frac{1}{P(x)} \quad (3.1)$$

where the set  $S$  spans the entire space of possible discrete values of  $x$ . We can go on to derive upper and lower bounds for the entropy. Notice that  $\mathbf{H} \geq 0$  since  $P(x) \leq 1$  and therefore  $\log P(x) \leq 0$  for all  $x$ . At the same time, if we define a variable  $Y = \frac{1}{\log x}$ , we can write

$$\begin{aligned} \mathbf{H} &= \mathbf{E}[\log Y] \\ &\leq \log \mathbf{E}[Y] \\ &= \log \sum_y P(x) \frac{1}{P(x)} \\ &= \log |S| \end{aligned}$$

which is just the entropy of a uniform distribution.

### 3.2.1 Joint and Conditional Entropy

In this section, we discuss joint and conditional entropy which are really just two sides of the same coin

$$\begin{aligned} \mathbf{H}(X, Y) &= \sum_{x,y} P(x, y) \log \frac{1}{P(x, y)} \\ &= \sum_{x,y} P(x)P(y|x) \log \frac{1}{P(x)P(y|x)} \\ &= \sum_{x,y} P(x)P(y|x) \log \frac{1}{P(x)} + \sum_{x,y} P(x)P(y|x) \log \frac{1}{P(y|x)} \\ &= \sum_{x,y} P(x)P(y|x) \log \frac{1}{P(x)} + \sum_x P(x) \sum_y P(y|x) \log \frac{1}{P(y|x)} \\ &= H(X) + H(Y|X) \end{aligned}$$

This result defines the **chain rule** for entropy. We typically refer to the term  $H(Y|X)$  as the **conditional entropy**. It can be calculated independently using the following definition

$$\begin{aligned} H(X|Y) &= \mathbf{E}_y H(X|Y = y) \\ &= \mathbf{E}_y \sum_x P(X|Y = y) \log \frac{1}{P(X|Y = y)} \end{aligned}$$

Furthermore, it can be shown that the chain rule derived above applies to a tuple of random variables longer than two.

$$H(X_1, \dots, X_m) = H(X_1) + H(X_2|X_1) + H(X_3|X_2, X_1) \dots H(X_m|X_1 \dots X_{m-1})$$

Recalling that conditioning reduces entropy or does nothing at all, we can write

$$H(X_1, \dots, X_m) \leq H(X_1) + H(X_2) + \dots + H(X_m)$$

which is referred to as the **subadditivity** property of entropy. We should also address what to do when we need to compute the entropy of a joint distribution  $(X, Y)$  conditioned on a variable  $Z$  or when  $Z$  itself is conditioned on a joint distribution. These two things are related by using the chain rule for joint entropy

$$H(X, Y|Z) = H(X, Y) + H(Z|X, Y)$$

Now we will prove that conditioning the distribution of a random variable  $X$  on another variable  $Y$  i.e. can reduce the entropy of  $X$ . What we need to show is that  $H(X|Y) - H(X) \leq 0$ .

### 3.3 KL-Divergence and Mutual Information

The Kullback-Leiber distance or **KL Divergence** is a measure of the distance between two distributions over a random variable  $X$ . Assume we have two distributions  $P, Q$  on a random variable  $X$  where  $P$  is the correct distribution on  $X$  and  $Q$  is an incorrect distribution. By definition, the KL-Divergence  $D_{KL}(P||Q)$  is the extra information (bits) it takes to communicate  $X$  when using the incorrect distribution  $Q$ . To be precise,  $H(Q) = H(P) + D_{KL}(P||Q)$ .

**Definition 1.** *The KL-Divergence is*

$$D_{KL}(P||Q) = \sum_X P(X) \log \frac{P(X)}{Q(X)}$$

Additionally, we can show that KL-Divergence follows a chain rule

Furthermore, an indispensable tool in information theory is the idea of **mutual information** which, as the name suggests, measures the amount of overlapping information in a pair of random variables. More formally, it is the KL-Divergence between the joint distribution of the pair of variables and the product of their marginal distributions (which implies they are independent)

**Definition 2.** *The mutual information is*

$$\begin{aligned} I(X;Y) &= D_{KL}(P(X,Y)||P(X)P(Y)) \\ &= \sum_x \sum_y P(X,Y) \log \frac{P(X,Y)}{P(X)P(Y)} \end{aligned}$$

A very useful property of the mutual information is that it is strongly related to conditional entropy and statistical independence. Conditional entropy tells us how much information is contained in a variable  $X$  which its distribution is conditioned on  $Y$ . We might expect that this conditioning doesn't really have an effect if  $X$  and  $Y$  are completely independent. Indeed,

$$\begin{aligned} I(X;Y) &= D_{KL}(P(X,Y)||P(X)P(Y)) \\ &= \sum P(x,y) \log \frac{P(x,y)}{P(x)P(y)} \\ &= \sum P(x) \log \frac{1}{P(x)} + \sum P(x,y) \log \frac{P(x,y)}{P(y)} \\ &= H(X) - H(X|Y) \end{aligned}$$

Note that this result implies that  $I(X;Y) = I(Y;X)$ . We will next address the mutual information between a distribution on  $X$  and a joint distribution  $(Y, Z)$  making use of the relationship derived above.

$$\begin{aligned} I(X;(Y,Z)) &= H(X) - H(X|Y,Z) \\ &= H(X) + H(Y,Z|X) - H(Y,Z) \end{aligned}$$

Finally, we look at the mutual information between a distribution on  $X$  and a conditional distribution  $Y|Z$ .

### 3.3.1 The Data-Processing Inequality

The data-processing inequality states that if a function operates on a random variable  $X$  it can only decrease its entropy. That is, for any function  $f$  s.t.  $Y = f(X)$ , we have that  $H(Y) \geq H(X)$ . We can prove that this is true using the mutual information  $I(X;Y)$ .

## 3.4 Source Coding

**Definition 3.** *A code of a set  $S$  that uses an alphabet  $\Omega$  is a map  $C : S \rightarrow \Omega$  that assigns each element of  $S$  a finite string over the alphabet  $\Omega$ . We say that the mapping  $C$  is **prefix free** if for all pairs  $x, y \in S$  where  $x \neq y$ ,  $C(x)$  is not a prefix of  $C(y)$ .*

Most of the time the alphabet  $\Omega$  we use is the set  $0, 1$ .

### 3.4.1 Kraft's Inequality

**Definition 4.** For a binary code, there exists a prefix free code  $C$  with codeword lengths  $l_i$  if and only if

$$\sum_i 2^{-l_i} \leq 1 \quad (3.2)$$

At this point we would like to apply the concept of entropy to source coding. Indeed, it is true that if we have a random variable  $X$  over the set  $S$ , the minimum number of bits it will take us to communicate the value of  $X$  on average is the entropy  $H(X)$ .

*Proof.* The expected number of bits to communicate  $X$  is given by  $\sum_x p(x)|C(x)|$

$$\begin{aligned} H(X) - \sum_x P(x)|C(x)| &= \sum_x P(x) \left[ \log \frac{1}{P(x)} - |C(x)| \right] \\ &= \sum_x P(x) \log \frac{1}{P(x) 2^{|C(x)|}} \\ &\geq \log \sum_x P(x) \frac{1}{P(x) 2^{|C(x)|}} \\ &= \log \sum_x \frac{1}{2^{|C(x)|}} \\ &\leq 0 \end{aligned} \quad \square$$

by Kraft's inequality for prefix-free codes.

### 3.4.2 Source Coding Theorem

So far we have seen how to construct a prefix-free code and that the absolute lower bound on the number of bits it takes to encode a random variable is its entropy. Next, we would like to answer the following question: how do we actually design a code to communicate a random variable  $X$  so that it approaches this lower bound? The answer is addressed by the *fundamental source coding theorem*

**Theorem 1.** For all  $\epsilon > 0$  there is a  $n_0 \leq n$  such that given  $n$  instances of a variable  $X$  it is possible to communicate  $X$  with  $H(X) + \epsilon$  bits on average.

This means that we can approach the entropy by increasing  $n$ .

## 3.5 Error Correcting Codes (Channel Coding)

Suppose that we have a message  $W$  that we would like to communicate. Preceding communication,  $W$  must be encoded according to a well-defined

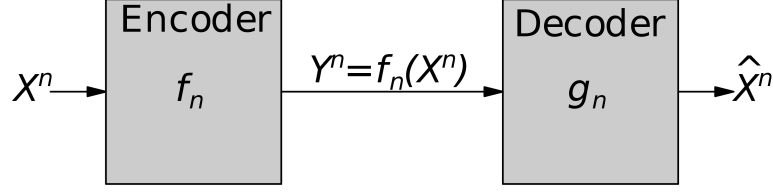


Figure 3.1: Channel coding Markov chain

encoding scheme. In other words we need to define a mapping  $W \rightarrow X^n$  which could be a sequence of  $n$  bits. The bit string  $X^n$  is then communicated over a channel and, at the destination, decoded using an inverse mapping  $Y^n \rightarrow \hat{W}$ . The challenge of channel coding is that, due to the fact that the channel can be noisy, it is possible that  $X^n \neq Y^n$  and therefore  $W \neq \hat{W}$ . Our primary goal is to determine optimal codes under the assumption that the channel can introduce a bounded error into the original message.

Lets say that we want to send one symbol of a family of symbols over a channel. The distribution over symbols has some entropy and therefore we can specify the minimum number of bits I must send for you to know which symbol I sent, on average. However, in a lossy channel, the entropy is difficult to achieve and so we define the *rate* which is the ratio of the number of bits it takes to specify the message (the entropy of a uniform distribution over messages) and the number of bits  $n$  I actually send you. We can imagine a much lower rate when the channel is very noisy.

$$R = \frac{\log M}{n}$$

In practice, we try to maximize the value of the rate. We can think of the communication channel as a conditional distribution  $P(Y|X)$ . At the same time, if we have a distribution over the inputs  $P(X)$  We define the *channel capacity* as the following

$$C = \max_{P(X)} I(X; Y)$$

If you recall that  $I(X; Y) = H(X) - H(X|Y)$  which is zero when knowledge of  $Y$  affords no information about  $X$  and is at most the entropy of  $X$  (when  $Y = X$ ). Of course, the channel is useless if  $H(X|Y) = H(X)$ . The channel coding theorem says that the maximum achievable rate is the channel capacity:

$$\sup R = \max_{P(X)} I(X; Y)$$



## 4

# The Neural Code

## 4.1 Introduction

Throughout this section of the text we will coarsely examine the biophysical properties of the elements of the nervous systems: neurons. The hope is that by introducing the nervous system from a physical perspective, much of the mathematical developments in later chapters will seem natural and well-motivated. However, we will see that much of the later work neglects many of the features of neurons this chapter will introduce and for good reason. It never hurts to know what you have left out of your model.

Neurons are one of the most interesting cell types in animals particularly for their evolved ability to sense their external environment and produce a myriad of responses in favor of the survival of the host organism. All vertebrates have a distinct central nervous system consisting of a brain and a spinal cord which process information from or send information to a peripheral nervous system. Here we will discuss the cellular elements that allow neurons to send information to each other in such a system.

## 4.2 Neuroelectronics

Neurons are bounded by a 3-4nm thick membrane composed of a diverse set of lipids forming a bi-layer and are submerged in a bath of water and ions *in vivo*. These cells have evolved a vast array of membrane ion channels that control the concentrations of ions on either side of the membrane. The degree to which these channels allow ions to flow in and out of the membrane is the basis of all neural processing and in turn all of our actions and conscious awareness.

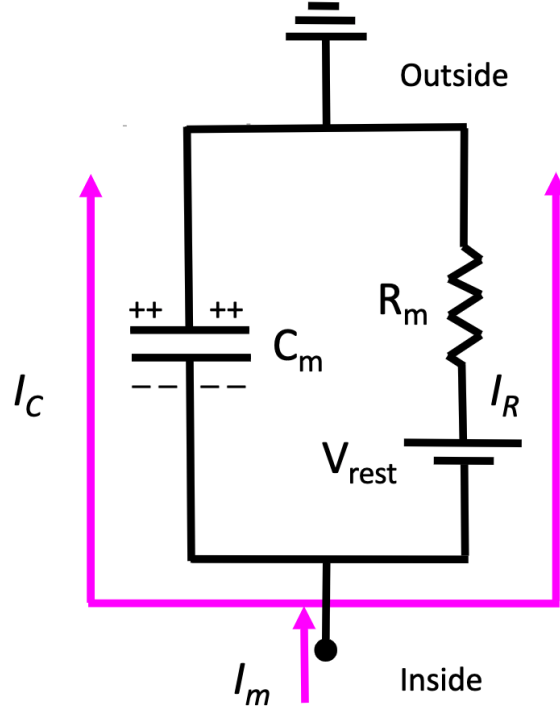


Figure 4.1: RC Circuit Model

#### 4.2.1 RC Model

A well-known model of a neuron membrane is an RC circuit where the resistor represents a sort of permeability of the membrane to ions while the capacitor models the membrane itself. Since either side of the membrane contains electrically charged ions there exists an electrical potential in these regions and in turn a potential difference across the lipid bi-layer. The current through the membrane is ultimately determined by the balance between diffusion of ions down their concentration gradients and electrical potential gradients.

If we inject a current  $I_E$  into the cell, by conservation of current we have

$$I_E = C \frac{dV}{dt} + I_R$$

where  $V$  is the membrane voltage. We can also use Maxwell's equations (Kirchhoff's voltage rule) to write

$$I_R = \frac{V - V_0}{R}$$

where we have replaced the resistance by the conductance  $g$ . Combining these two equations we are left with the following differential equation for the membrane voltage.

$$C \frac{dV}{dt} = I_E - \frac{V - V_0}{R}$$

which is commonly rewritten by defining the membrane time constant  $\tau = RC$

$$\tau \frac{dV}{dt} = I_E R - (V - V_0)$$

Notice that as  $t \rightarrow \infty$  we have  $V_C - V_{rest} \rightarrow I_E R + V_0$  and so we define the steady-state voltage  $V_\infty = I_E R + V_0$ . Inserting this new definition into the equation

$$\tau \frac{dV}{dt} = V_\infty - V$$

which can be solved to give

$$V(t) = V_\infty (1 - e^{-\frac{t}{\tau}})$$

#### 4.2.2 Equilibrium and Reversal Potentials

The flux of ions through the cell membrane is due to an imbalance of electrical and chemical potentials. Recall the Boltzmann distribution of energies at a temperature  $T$

$$P(\epsilon) = \frac{e^{-\frac{\epsilon}{k_B T}}}{Z}$$

where the partition function  $Z = \int e^{-\frac{\epsilon}{k_B T}} d\epsilon$ . We define the **reversal potential** as the membrane potential at which electrical forces perfectly balance the diffusion of ions down their concentration gradient and there is no current through the membrane. We compute this using the famous Nernst equation

$$E = \frac{RT}{zF} \ln \frac{[o]}{[i]}$$

Perhaps we would rather compute the resting potential of the membrane with knowledge of the concentrations of ions inside and outside of the cell. For a set of positive ions  $M$  and negative ions  $A$

$$E = \frac{RT}{F} \ln \frac{\sum_i M_i [M_i]_{out} + \sum_i A_i [A_i]_{in}}{\sum_i M_i [M_i]_{in} + \sum_i A_i [A_i]_{out}}$$

where we have used  $M_i, A_i$  to represent the permeability of these ions.

### 4.2.3 Membrane Current

We saw in our basic RC circuit model of the cell membrane that the membrane current  $I_R$  was very simply given by Ohm's law  $I_R = (V - V_0)/R$ . This very simple model assumes that the entire membrane can be lumped as a single resistor  $R$ ; however, in reality, neural membranes have a diverse set of channel types each with their own respective resistance to current flow. Therefore, we need to expand  $R$  into a set of parallel resistors each with their own resistance.

$$I = \sum_i g_i (V - E_i)$$

This setup is more accurate but still not entirely correct; Ohm's law assumes that we are dealing with a **linear circuit** where the resistance is a constant (not a function of the voltage). In the cell, many channels are **voltage gated** which means that they *are* a function of the voltage and Ohm's law doesn't apply. A more correct definition of the total membrane current is then

$$I = \sum_i g_i(V) (V - E_i) \quad (4.1)$$

The functions  $g_i(V)$  are difficult to estimate due to the complexity of ion channels and their mechanisms of opening and closing. Hodgkin and Huxley came up with a model for these voltage-dependent conductance for sodium and potassium. Their model defines what we call an **open probability** for a single ion channel at a particular voltage. In the case of the potassium channel, Hodgkin and Huxley thought of the channel as a structure composed of 4 particles which each probability  $n$  of being in the *permissive state*. If all four particles are in such a state, the channel will open. Each particle has a probability of being in this state which increases with the voltage sigmoidally

$$n(V) = \frac{1}{1 + (A_\beta/A_\alpha) \exp((B_\alpha - B_\beta)V/V_T)} \quad (4.2)$$

and so the total probability the channel is open is  $n^4$ . The model of the sodium channel is slightly different. They introduced two probabilities

$m$  and  $h$  where  $m$  is directly analogous to  $n$  and  $h$  is the probability that a gate on the sodium channel is open. They fit their experimental data using  $m^3$  and so probability the sodium channel is open is then  $m^3h$ . The total membrane current according to the Hodgkin-Huxley model can then be found by expanding Eq. 1.1

$$I = \bar{g}_{Na}m^3h(V - E_{Na}) + \bar{g}_Kn^4(V - E_K) + \bar{g}_L(V - E_L) \quad (4.3)$$

where a leakage conductance  $g_L$  is also included.