

# Information Theory & the Efficient Coding Hypothesis

Jonathan Pillow

Mathematical Tools for Neuroscience (NEU 314)  
Spring, 2016

lecture 19

# Information Theory

*A mathematical theory of communication,*  
Claude Shannon 1948

- Entropy
- Conditional Entropy
- Mutual Information
- Data Processing Inequality
- Efficient Coding Hypothesis (Barlow 1961)

# Entropy

$$H(x) = - \sum_x p(x) \log p(x)$$

averaged over  $p(x)$       “surprise” of  $x$

- average “surprise” of viewing a sample from  $p(x)$
- number of “yes/no” questions needed to identify  $x$  (on average)

for distribution on  $K$  bins,

- maximum entropy =  $\log K$  (achieved by uniform dist)
- minimum entropy = 0 (achieved by all probability in 1 bin)

# Entropy

$$\begin{aligned} H(x) &= - \sum_x p(x) \log p(x) \\ &= -\mathbb{E}[\log p(x)] \end{aligned}$$

# aside: log-likelihood and entropy

model:  $P(x|\theta)$

entropy H:  $-\mathbb{E}[\log P(x|\theta)]$

How would we compute a Monte Carlo estimate of this?

draw samples:  $x_i \sim P(x|\theta)$  for  $i = 1, \dots, N$

compute average:  $\hat{H} = -\frac{1}{N} \sum_{i=1}^N \log P(x_i|\theta)$

  
log-likelihood

- Neg Log likelihood = Monte Carlo estimate for entropy!
- maximizing likelihood  $\Rightarrow$  minimizing entropy of  $P(x|\theta)$

# Conditional Entropy

$$H(x|y) = - \sum_y p(y) \sum_x p(x|y) \log p(x|y)$$

averaged  
over  $p(y)$

entropy of  $x$  given  
some fixed value of  $y$

# Conditional Entropy

$$H(x|y) = - \sum_y p(y) \sum_x p(x|y) \log p(x|y)$$

averaged  
over  $p(y)$

entropy of  $x$  given  
some fixed value of  $y$

$$= - \sum_{x,y} p(x,y) \log p(x|y)$$

$$= H(x) \quad \text{if} \quad P(x,y) = P(x)P(y)$$

“On average, how uncertain are you about  $x$  if you know  $y$ ?”

# Mutual Information

$$I(x, y) = H(x) - H(x|y)$$

total entropy in X minus  
conditional entropy of X given Y

$$= H(y) - H(y|x)$$

total entropy in Y minus  
conditional entropy of Y given X

$$= H(x) + H(y) - H(x, y)$$

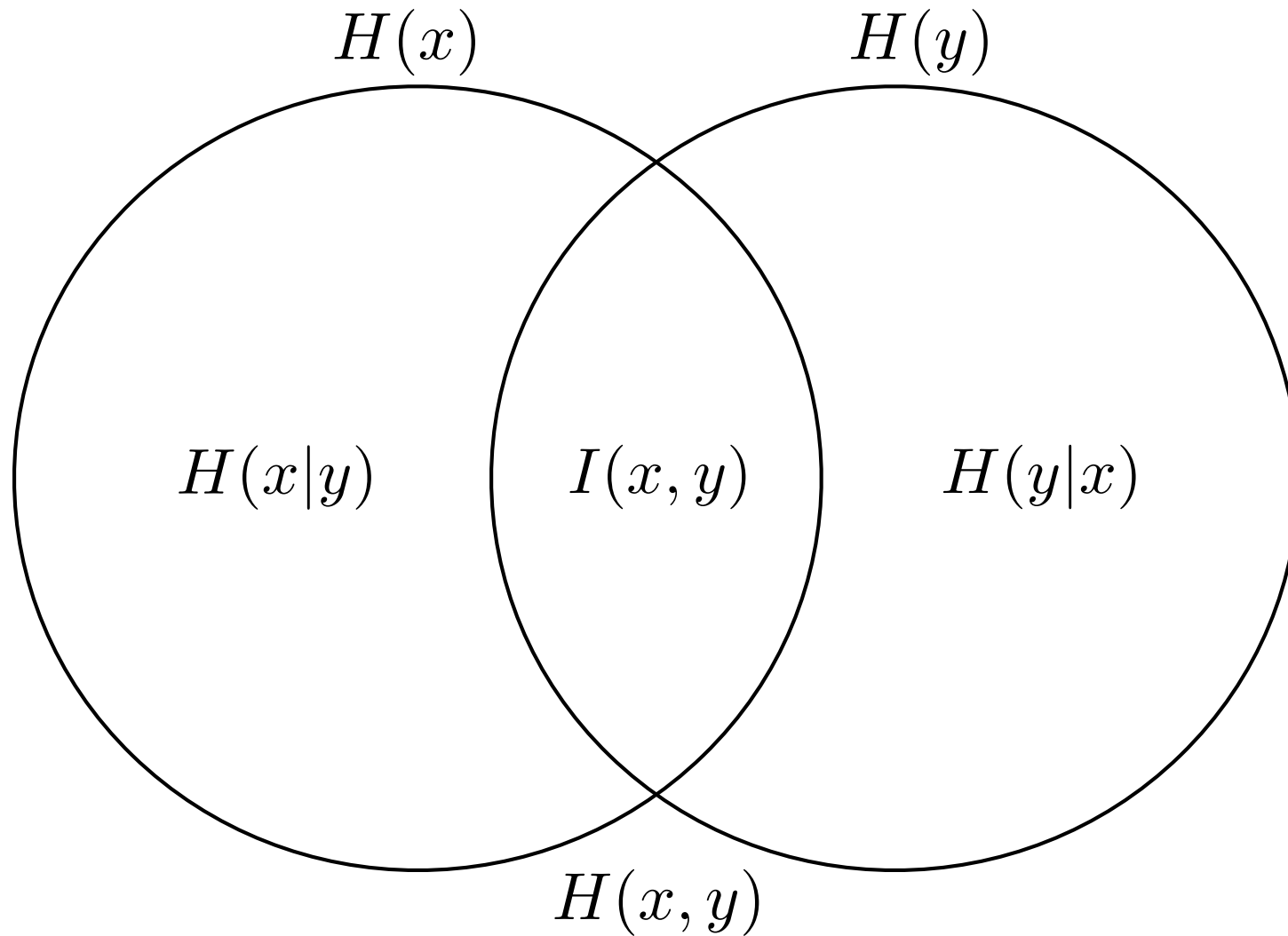
sum of entropies  
minus joint entropy

“How much does X tell me about Y (or vice versa)?”

“How much is your uncertainty about X reduced from knowing Y?”



# Venn diagram of entropy and information



# Data Processing Inequality

Suppose  $S \rightarrow R_1 \rightarrow R_2$  form a Markov chain, that is

$$P(R_1, R_2|S) = P(R_2|R_1)P(R_1|S)$$

Then necessarily:  $I(S, R_2) \leq I(S, R_1)$

- in other words, we can only lose information during processing

# Efficient Coding Hypothesis:

- goal of nervous system: maximize information about environment  
(one of the core “big ideas” in theoretical neuroscience)

**redundancy:**  $R = 1 - \frac{I}{C}$

↖ mutual information  
↖ channel capacity

# Efficient Coding Hypothesis:

- goal of nervous system: maximize information about environment  
(one of the core “big ideas” in theoretical neuroscience)

**redundancy:**  $R = 1 - \frac{I}{C}$

↖ mutual information  
↖ channel capacity

**mutual information:**

$$I(x, y) = H(y) - H(y|x)$$

response entropy    “noise” entropy

- avg # yes/no questions you can answer about x given y (“bits”)

**channel capacity:**

$$C = \sup_{P_x} I(x, y)$$

- upper bound on mutual information
- determined by physical properties of encoder

# Barlow's original version:

**redundancy:**  $R = 1 - \frac{I}{C}$  ↖ mutual information

**mutual information:**

$I(x, y) = H(y) - \cancel{H(y|x)}$  if responses are noiseless

response entropy    "noise" entropy

# Barlow's original version:

**redundancy:**  $R = 1 - \frac{H(Y)}{C}$  ← response entropy

**mutual information:**

$$I(x, y) = H(y) - \cancel{H(y|x)} \quad \text{noiseless system}$$

response entropy    “noise” entropy

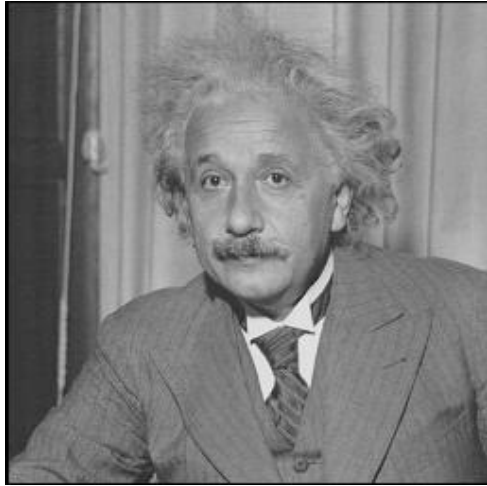
⇒ brain should maximize response entropy

- use full dynamic range
- decorrelate (“reduce redundancy”)

- mega impact: huge number of theory and experimental papers focused on decorrelation / information-maximizing codes in the brain

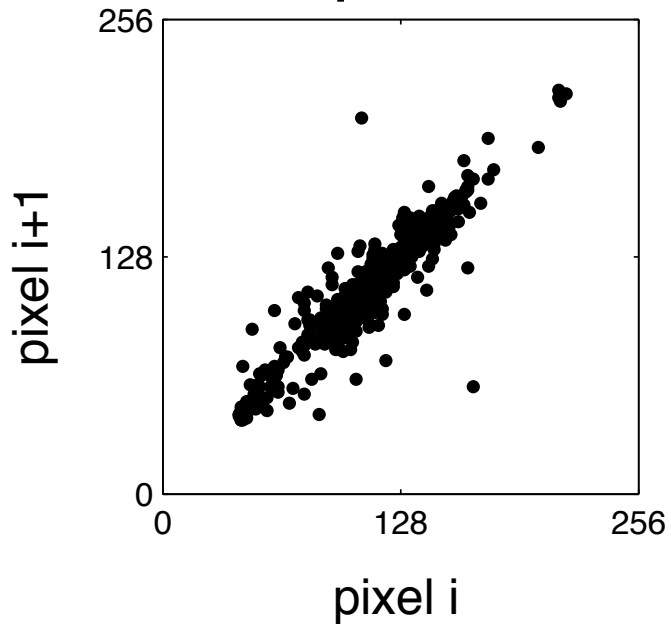
# basic intuition

natural image

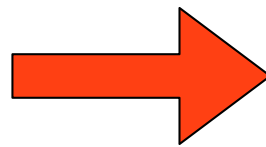


nearby pixels exhibit strong dependencies

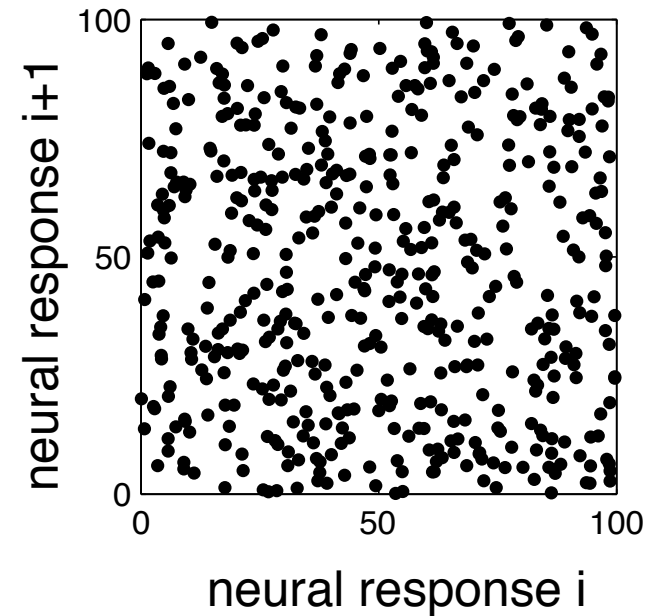
pixels



desired encoding



neural representation



Example: single neuron encoding stimuli from a distribution  $P(x)$

stimulus prior  $x \sim P(x)$

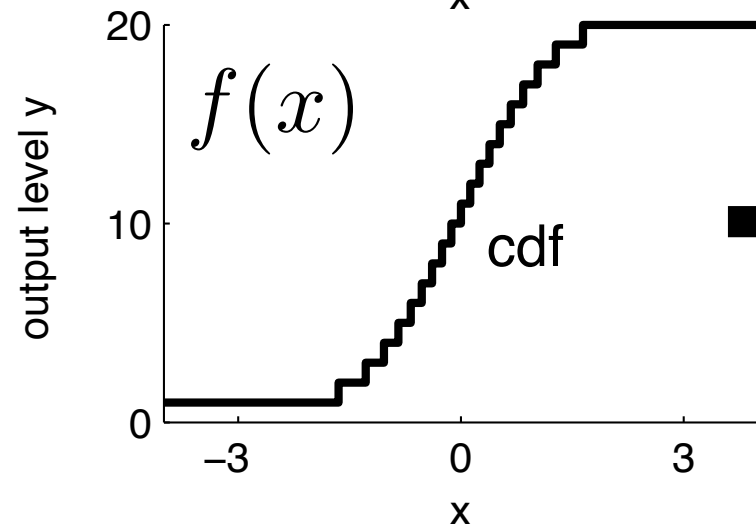
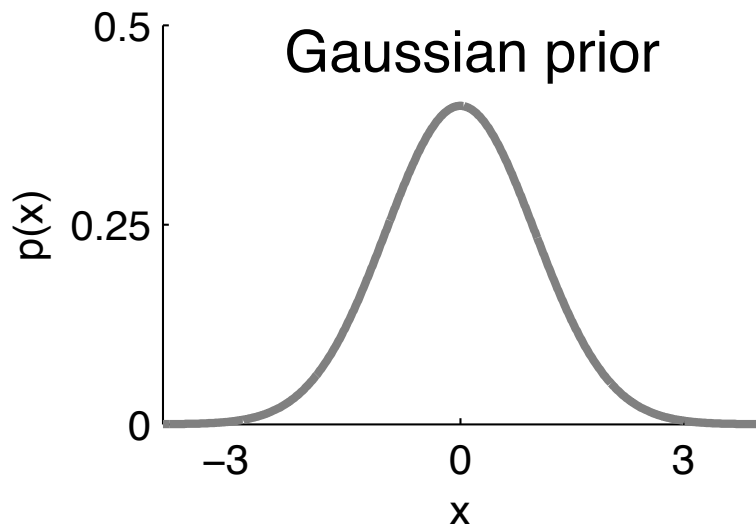
noiseless, discrete encoding  $y = f(x)$  (with constraint on range of  $y$  values)



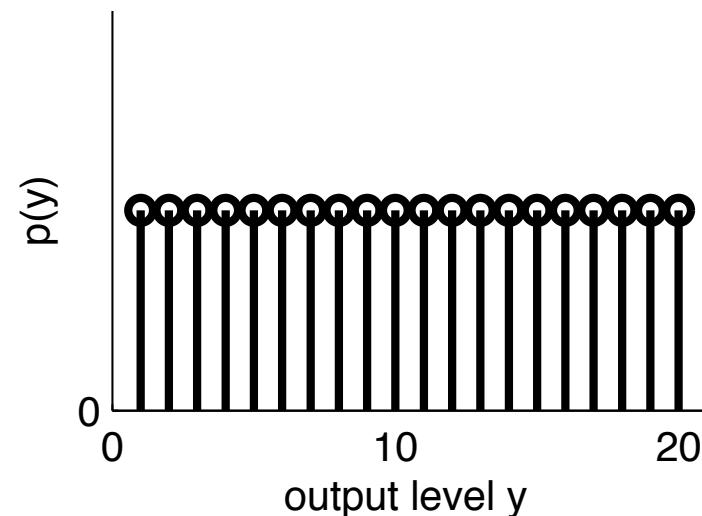
# Application Example: single neuron encoding stimuli from a distribution $P(x)$

stimulus prior  $x \sim P(x)$

noiseless, discrete encoding  $y = f(x)$  (with constraint on range of  $y$  values)

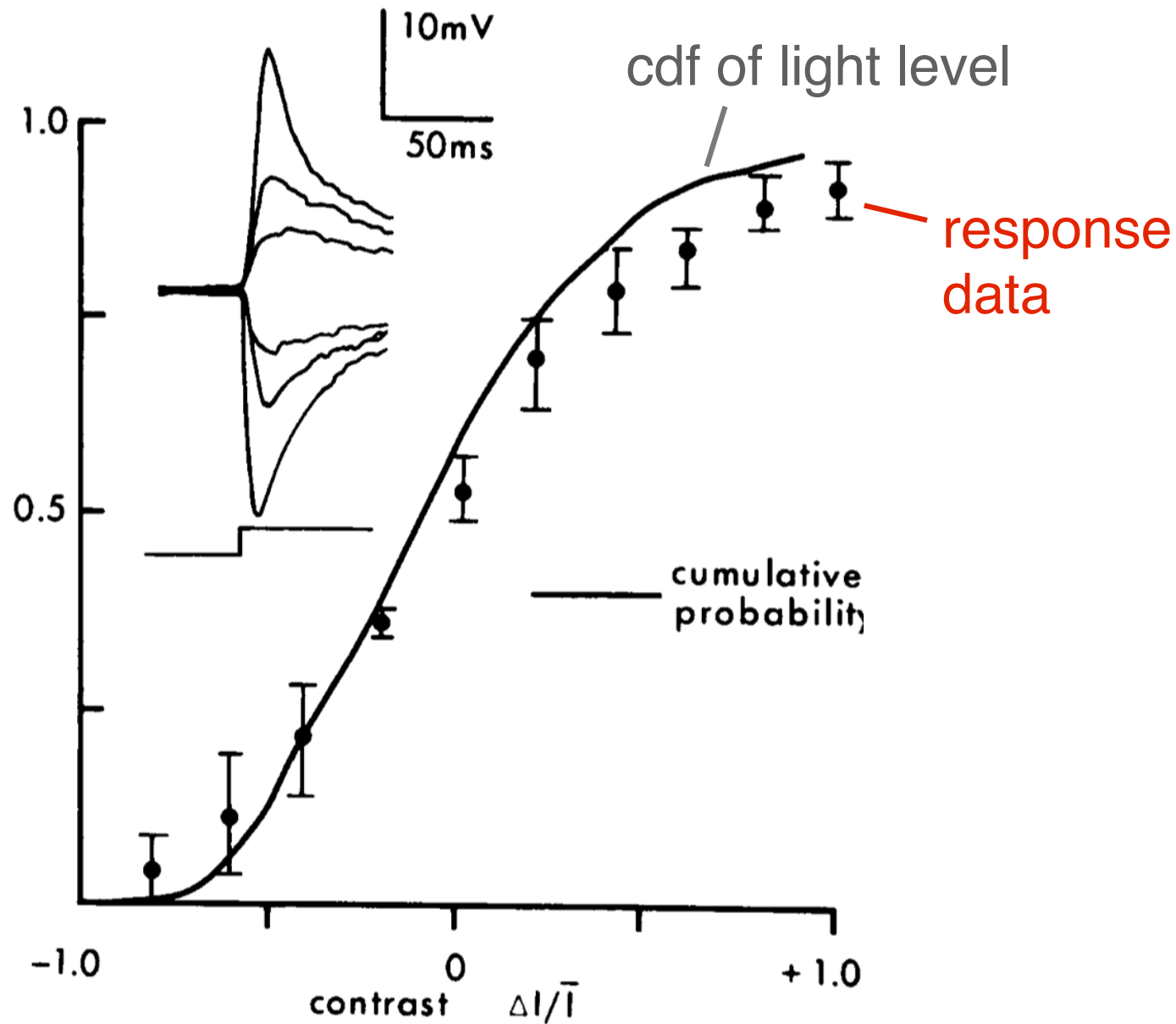


response distribution



# Laughlin 1981: blowfly light response

- first major validation of Barlow's theory



# summary

- entropy
- negative log-likelihood / N
- conditional entropy
- mutual information
- data processing inequality
- efficient coding hypothesis (Barlow)
  - neurons should “maximize their dynamic range”
  - multiple neurons: marginally independent responses
- direct method for estimating mutual information from data