

Prob modelling, deep learning for . . .

# ORIE 4742 - Info Theory and Bayesian ML

---

January 21, 2020

Semester: Spring 2020

# essential course information

## OH-TBA

- *instructor:* Sid Banerjee, [sbanerjee@cornell.edu](mailto:sbanerjee@cornell.edu)
- *TAs:* Cameron Ibrahim, [cai29@cornell.edu](mailto:cai29@cornell.edu)  
Shengyuan Hu, [sh797@cornell.edu](mailto:sh797@cornell.edu)
- *lectures:* TR 8:40-9:55am, Upson 216
- *website*  
<https://piazza.com/cornell/spring2020/orie4742>

## the fine print

- *grading*  
50% homeworks, 20% prelim, 25% project,  
5% class participation
- *homeworks*  
6 homeworks (on average 2 weeks for each)  
teams of up to 3  
submit single Jupyter notebook, with theory answers in Markdown  
typically due Monday 5pm, on <https://cmsx.cs.cornell.edu>  
4 late days across homeworks, lowest grade dropped
- *prelim*  
in class, tentatively March 26 (before spring break)  
no final exam
- *project*  
use techniques learned in class on ML problem of your choosing  
teams of up to 3, report due on final exam date

## what is this class about

- Q1. given data, how can we learn how it was generated?

## what is this class about

- Q1. given data, how can we learn how it was generated?
- Q2. how can we translate data and models into future decisions?

## what is this class about

- Q1. given data, how can we learn how it was generated?
- Q2. how can we translate data and models into future decisions?
- Q3. what are the fundamental limits and design principles of data-driven learning and decision-making

# what is this class about

- Q1. given data, how can we learn how it was generated?
- Q2. how can we translate data and models into future decisions?
- Q3. what are the fundamental limits and design principles of data-driven learning and decision-making

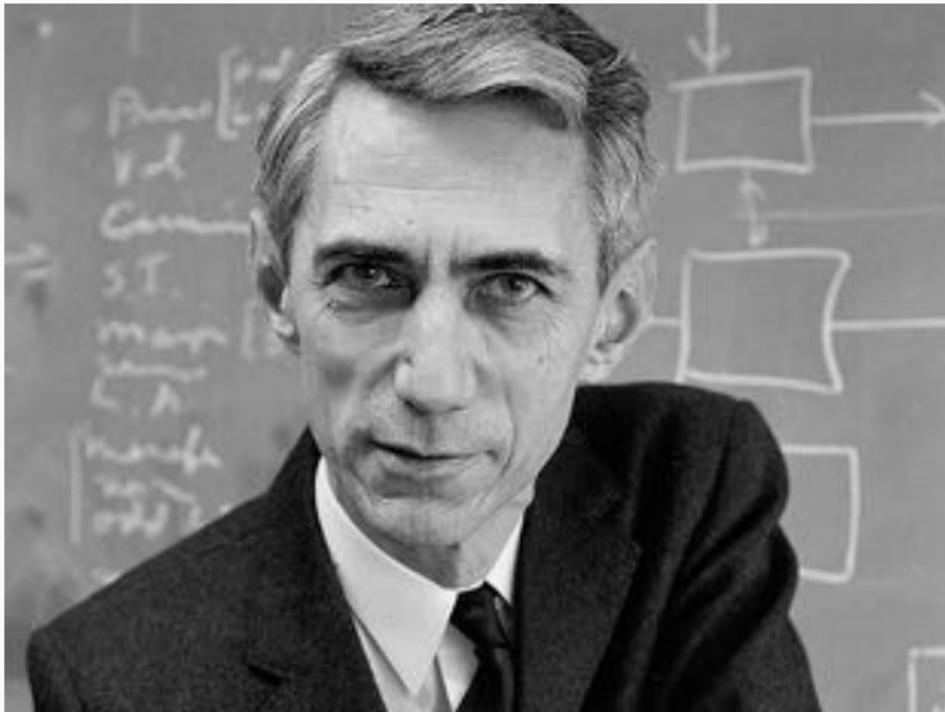
our approach in this course: **probabilistic modeling**

- bayesian inference: unified paradigm for learning and decision-making
- information theory: tool for designing and understanding data systems

model - encode 'inductive biases'  
- Priors on unknown quantities

data      update  
priors

**problem: communicating over a noisy channel**



**reading assignment: chapter 1 of Mackay**

# communicating over channels

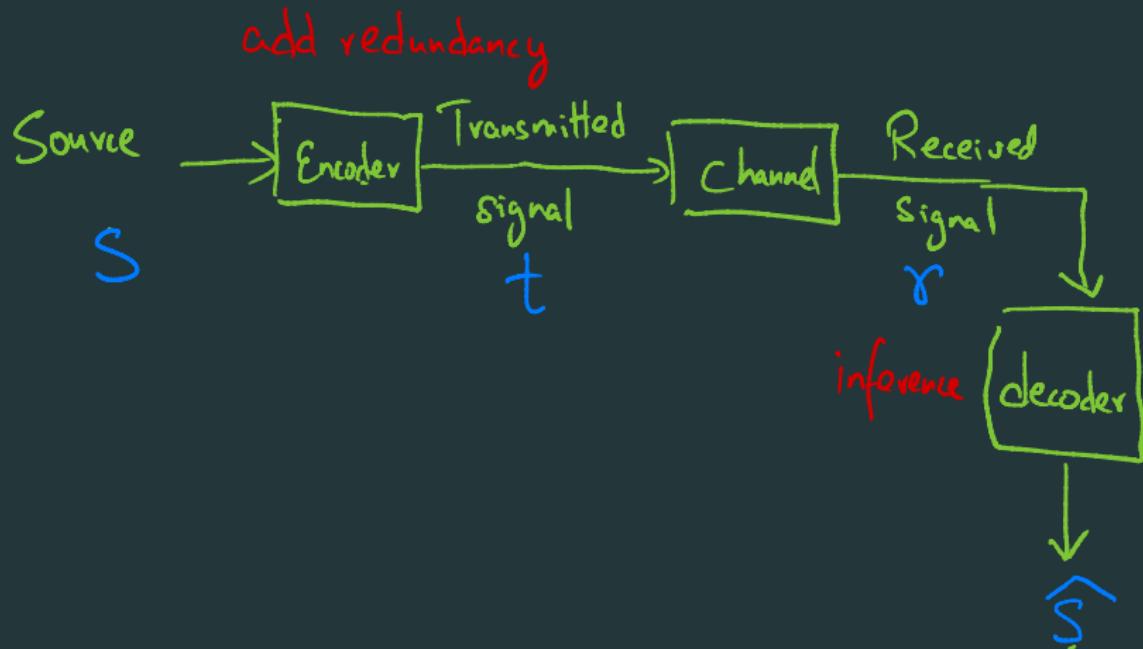
- mouth  $\xrightarrow{\text{air}}$  ear
- ear  $\xrightarrow[\dots]{\text{retina, nerves}}$  brain
- dna  $\xrightarrow{\text{reproduction}}$  dna
- cellphone  $\xleftarrow{\text{air}}$  basetower
- terrestrial antenna  $\xleftarrow{\text{space}}$  Mars rover
- data  $\xrightarrow{\text{storage}}$  data in future
- generative model  $\xrightarrow[\text{collection}]{\text{data}}$  data

$$\text{Signal} = \text{data} + \text{noise}$$

Two approaches

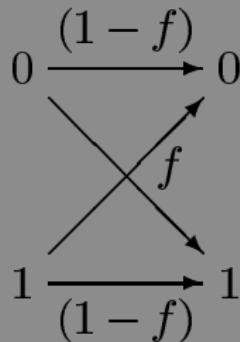
- physical approach
- systems approach

## the system's solution



## a toy model: the **binary symmetric channel**

$$f \equiv P[\text{bit flip}] \sim 0.1$$



$$\{0,1\}^n$$

$$P[r=0|t=0] = 1-f, P[r=1|t=0] = f$$

$$P[r=1|t=1] = f, P[r=0|t=1] = 1-f$$

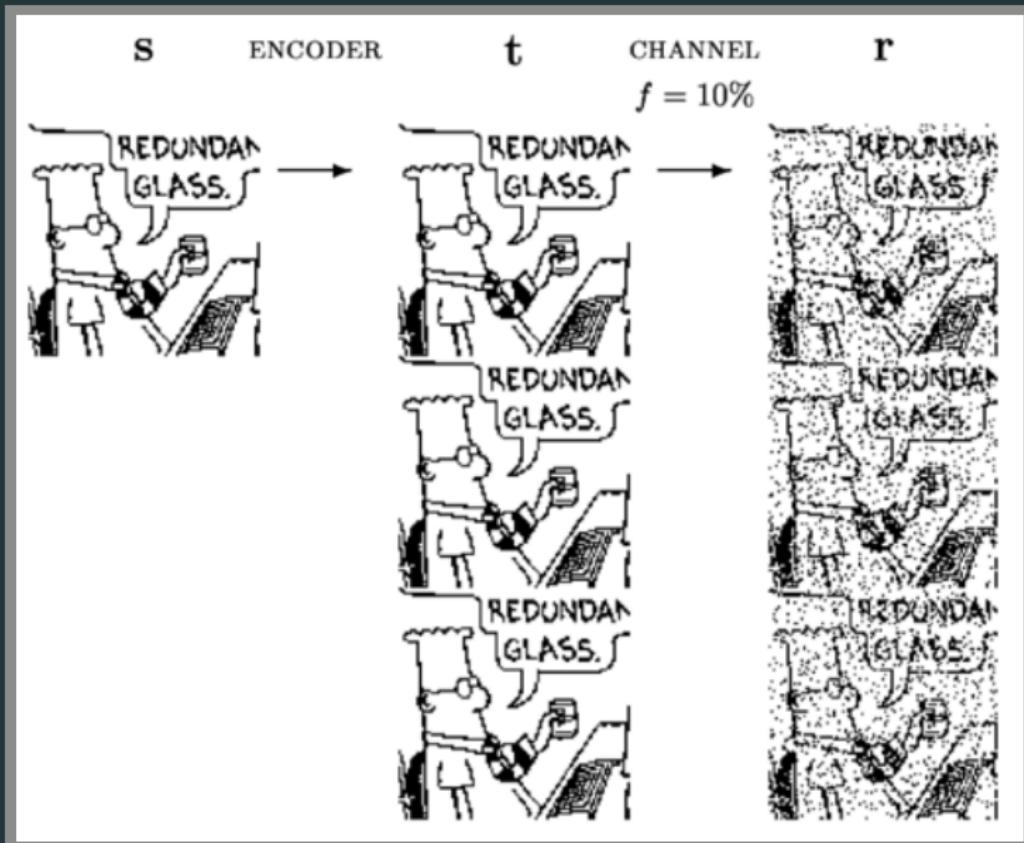
credit: David Mackay

## ideas for encoding

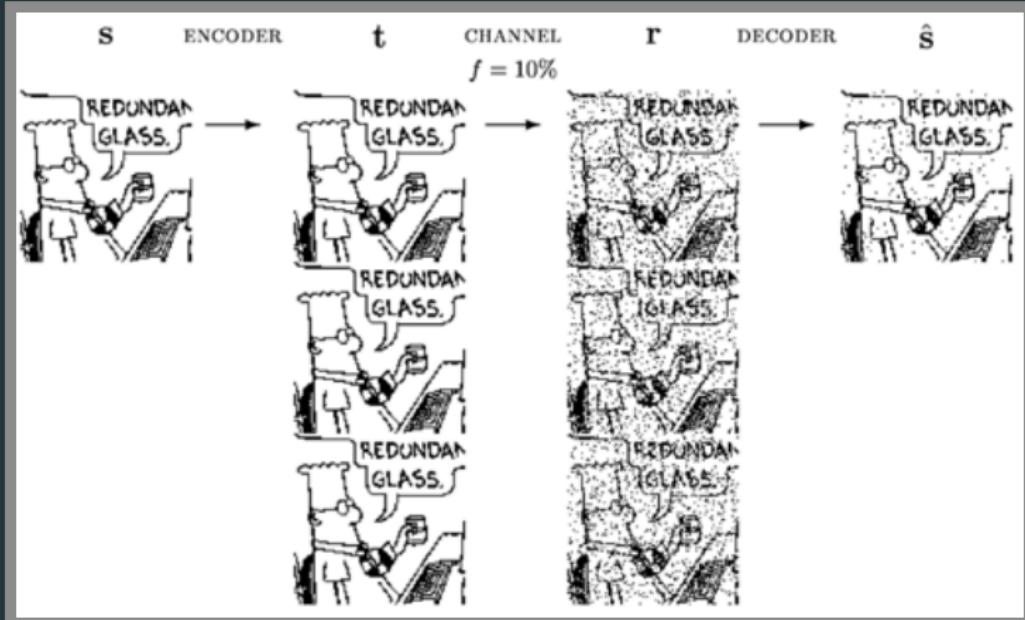
$$s = 011 \quad , \quad \text{noise} = 00100001100$$

- $t = s \oplus n \bmod 2$
- no encoding -  $r = 010$
- repeat 2 times -  $t = 00\ 11\ 11$ ,  $r = 00\ 01\ 11$   
3 times -  $t = 000\ 111\ 111$ ,  $r = 001\ 111\ 100$   
decoding - majority
- Parity bits -  $t_1 = s_1, t_2 = s_2, t_3 = s_3, t_4 = s_1 + s_2 + s_3 \bmod 2$   
 $t = 0110$ ,  $r = 0100$

## repetition codes: encoding



# repetition codes: decoding



credit: David Mackay

decoding rule = majority

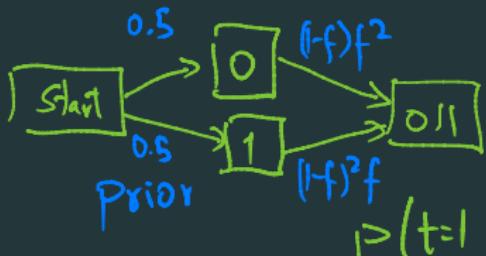
# repetition codes: inference (inverse probability)

- 3-repetitions,  $r = 011$

$$P(r=011 \mid t=0) = (1-f) \cdot f \cdot f = (1-f) f^2 \rightarrow \text{# of bit flips}$$

$$P(r=011 \mid t=1) = f \cdot (1-f) \cdot (1-f) = (1-f)^2 f \rightarrow \text{# of bit flips}$$

- Bayes thm •  $P(t=0 \mid r=011) = \frac{P(t=0) P(r=011 \mid t=0)}{\sum_{i=0,1} P(t=i) P(r=011 \mid t=i)}$



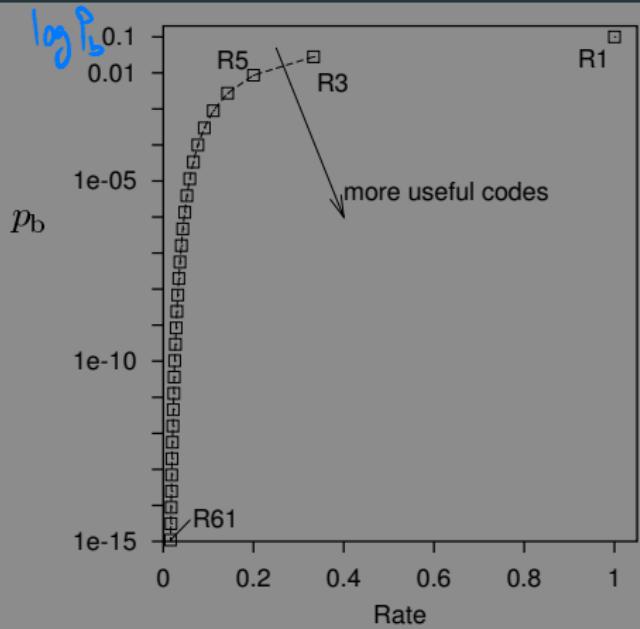
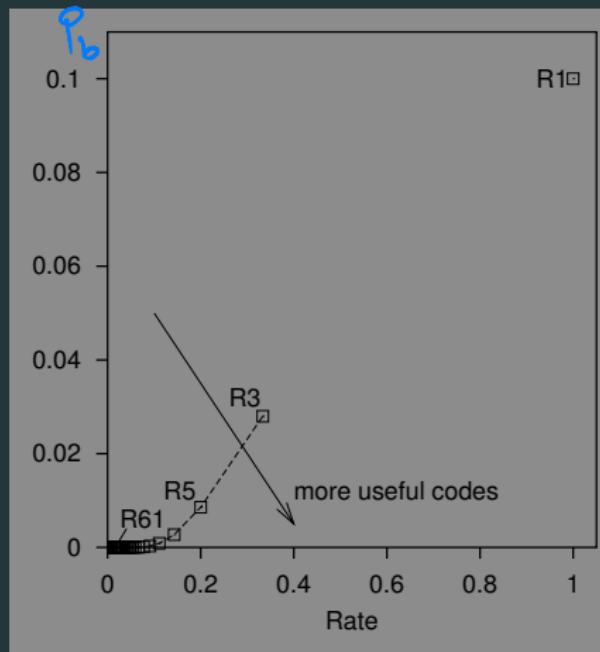
$$P(t=1 \mid r=011) = \frac{P(t=1) P(r=011 \mid t=1)}{Z}$$

- Inference  $\equiv$  choosing signal that best explains data normalization constant
- If priors equal, optimal inference rule  $\equiv$  majority

## repetition codes: performance

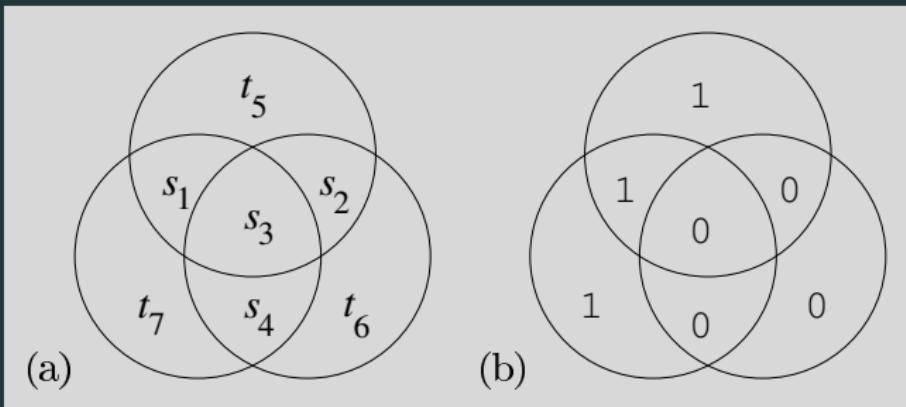
- If we use optimal inference rule, what is the probability of bit error (as a fn of # of reps R)
- $0 \xrightarrow{R=5} 00000 \xrightarrow{\text{BSC}} r \xrightarrow{\text{majority}} \hat{s}$ 
$$\begin{aligned}\mathbb{P}[\hat{s} \neq s] &= \mathbb{P}[\# \text{ of flips} \geq 3] \\ &= \binom{5}{3} f^3 (1-f)^2 + \binom{5}{4} f^4 (1-f) + f^5 \\ &\approx 10f^3.\end{aligned}$$
- In general  $\mathbb{P}[\hat{s} \neq s] \approx c f^{\lfloor R/2 \rfloor}$

## repetition codes: the rate-error plot



credit: David Mackay

# the (7,4) Hamming code



credit: David Mackay

## the (7,4) Hamming code: performance

s	t	s	t	s	t	s	t
0000	0000000	0100	0100110	1000	1000101	1100	1100011
0001	0001011	0101	0101101	1001	1001110	1101	1101000
0010	0010111	0110	0110001	1010	1010010	1110	1110100
0011	0011100	0111	0111010	1011	1011001	1111	1111111

credit: David Mackay

## the (7,4) Hamming code: performance

s	t	s	t	s	t	s	t
0000	0000000	0100	0100110	1000	1000101	1100	1100011
0001	0001011	0101	0101101	1001	1001110	1101	1101000
0010	0010111	0110	0110001	1010	1010010	1110	1110100
0011	0011100	0111	0111010	1011	1011001	1111	1111111

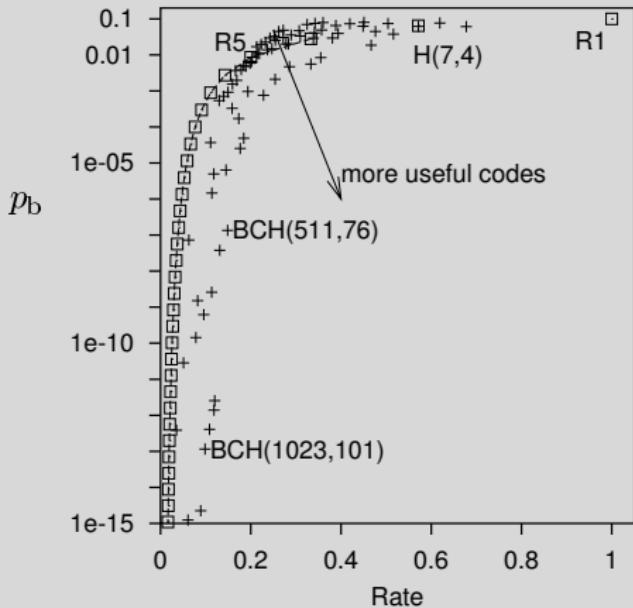
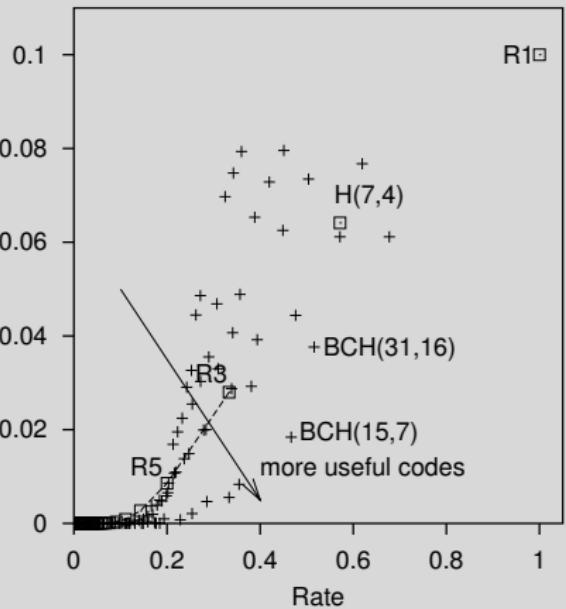
credit: David Mackay

### distance between codewords

the minimal Hamming distance between any two correct codewords is 3

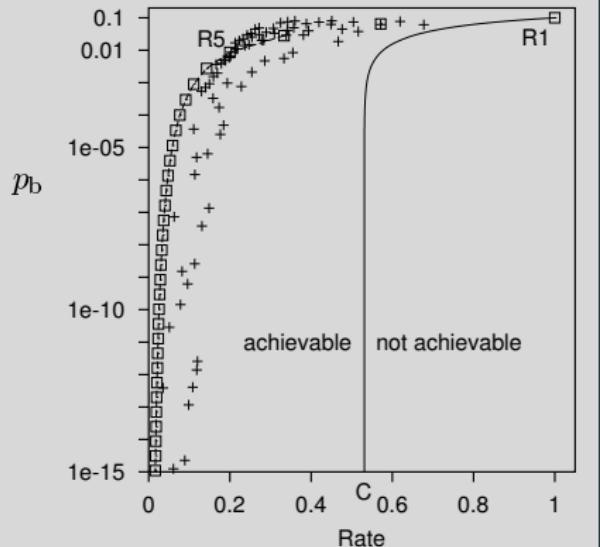
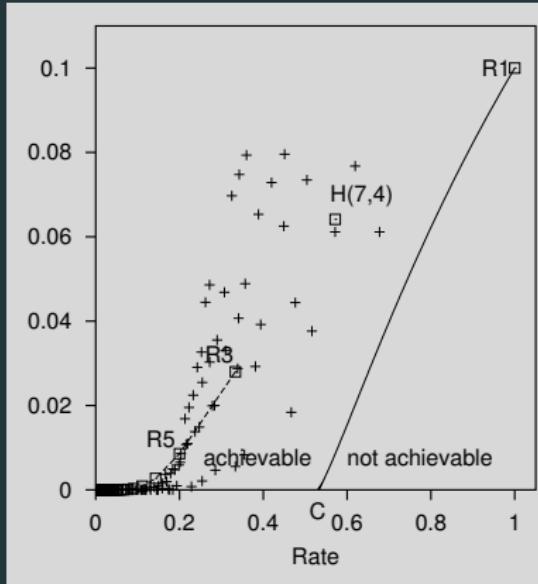
corollary - (7,4) Hamming code 'fixes' 1 error

## the rate-error plot



credit: David Mackay

# Shannon's channel coding theorem

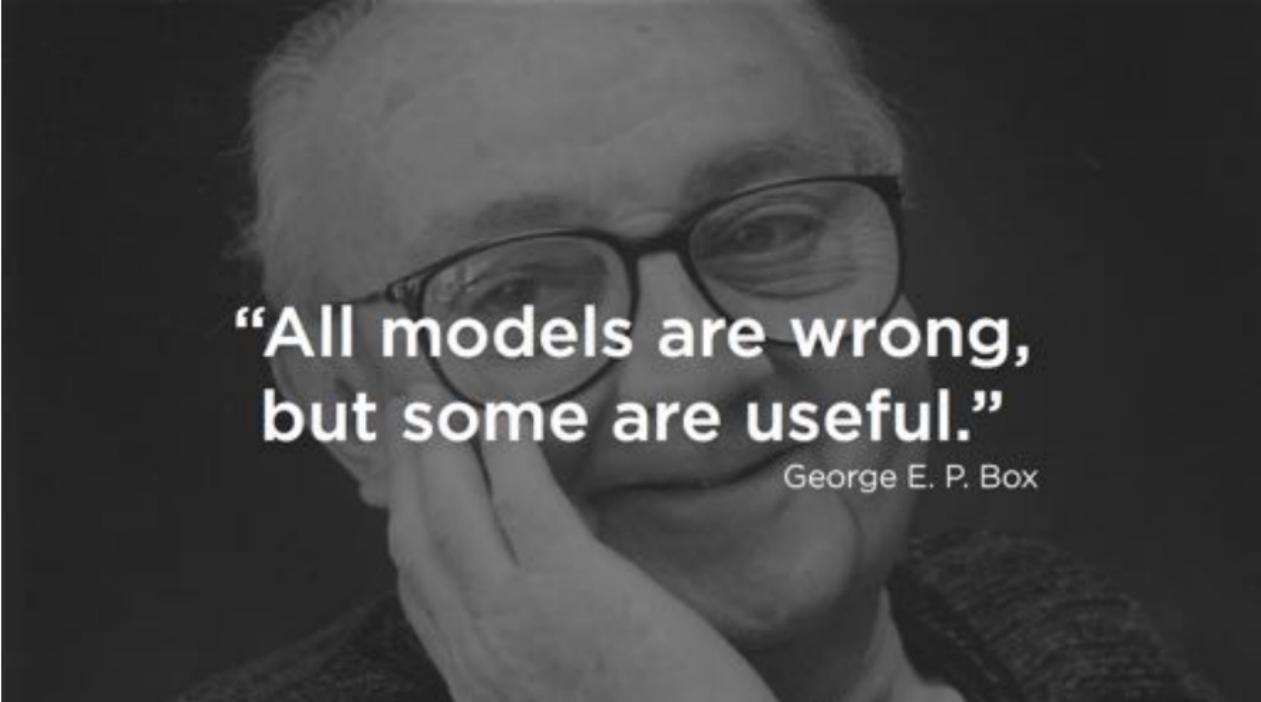


**Theorem (Claude Shannon, 1948)**

for any channel, 0-error communication is possible at a rate up to  $C > 0$

For BSC -  $C = H_2(f) = f \log_2\left(\frac{1}{f}\right) + (1-f) \log_2\left(\frac{1}{1-f}\right)$

# noisy channel communication $\leftrightarrow$ machine learning



**“All models are wrong,  
but some are useful.”**

George E. P. Box

## redundancy ⇒ inference

Emma Woodh\*use, hands\*me, clever\* and rich,\*with a  
comfortab\*e home an\* happy di\*position,\*seemed to\*unite som\*  
of the b\*st bless\*ngs of e\*istence;\*and had \*ived nearly  
twenty \*ne year\* in the\*world w\*th very\*little \*o distr\*ss  
or vex\*her. \*he was\*the yo\*ngest \*f the \*wo dau\*hters \*f a  
most \*ffect\*onate\* indu\*gent \*ather\* and \*ad, i\* cons\*qunc\*  
of h\*r si\*ter'\* mar\*iage\* bee\* mis\*ress\*of h\*s ho\*se f\*om a  
ver\* ea\*ly \*eri\*d. \*er \*oth\*r h\*d d\*ed \*oo \*ong\*ago\*for\*her  
to\*hate \*or\* t\*an\*an\*in\*i\*is\*int \*em\*mb\*an\*e \*f \*er\*ca\*es\*es\*  
a\*d\*h\*r\*p\*a\*e\*h\*d\*b\*e\* \*u\*p\*i\*d\*b\* \*n\*e\*c\*l\*e\*t\*w\*m\*n\*a\*  
g\*\*\*e\*\*\*,\*\*\*h\*\*\*h\*\*\* \*\*\*l\*\*\*n\*\*\*i\*\*\*l\*\*\*s\*\*\*r\*\*\*o\*\*\*a\*\*\*o\*\*\*e\*\*\*i\*  
a\*\*\*c\*\*\*n\*\*\*S\*\*\*e\*\*\*y\*\*\*s\*\*\*d\*\*\*s\*\*\*a\*\*\*r\*\*\*e\*\*\*n\*\*\*  
W\*\*\*o\*\*\*s\*\*\*i\*\*\*l\*\*\*a\*\*\*g\*\*\*n\*\*\*t\*\*\*a\*\*\*e\*\*\*v\*\*\*

credit: David Mackay

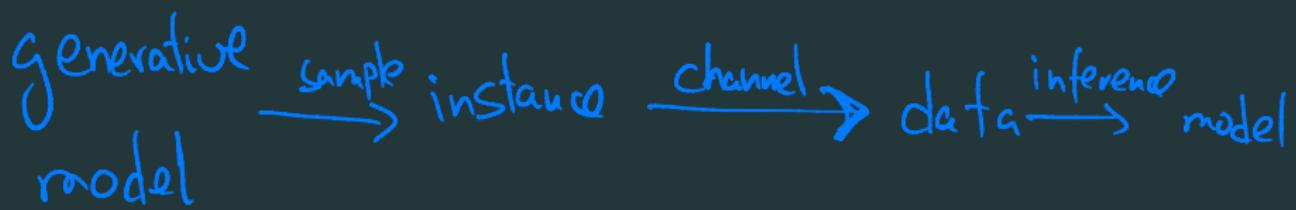
## redundancy ⇒ inference

Emma Woodh\*use, hands\*me, clever\* and rich,\*with a  
comfortab\*e home an\* happy di\*position,\*seemed to\*unite som\*  
of the b\*st bless\*ngs of e\*istence;\*and had \*ived nealy  
twenty \*ne year\* in the\*world w\*th very\*little \*o distr\*ss  
or vex\*her. \*he was\*the yo\*ngest \*f the \*wo dau\*hters \*f a  
most \*ffectionate\* indu\*gent \*ather\* and \*ad, i\* cons\*quenc\*  
of h\*r si\*ter\*\* mar\*riage\* bee\* mis\*ress\*of h\*s ho\*se f\*om a  
ver\* ea\*ly \*eri\*d. \*er \*oth\*r h\*d d\*ed \*oo \*ong\*ago\*for\*her  
to\*have \*orr t\*an\*an\*in\*is\*int \*em\*mb\*an\*e \*f \*er\*ca\*es\*es\*  
a\*d\*h\*r\*p\*a\*e\*h\*d\*b\*e\* \*u\*p\*i\*d\*b\* \*n\*e\*c\*l\*e\*t\*w\*m\*n\*a\*  
g\*\*\*e\*\*\*, \*\*h\*\*\*h\*\* \*l\*\*n\*\*i\*\*l\*\*s\*\*r\*\*o\*\*a\*\*o\*\*e\*\*\*i\*  
a\*\*\*c\*\*\*n\*\*\*S\*\*\*e\*\*\*y\*\*\*s\*\*\*d\*\*\*s\*\*\*a\*\*\*r\*\*\*e\*\*\*n\*\*\*  
W\*\*\*o\*\*\*s\*\*\*i\*\*\*l\*\*\*a\*\*\*g\*\*\*n\*\*\*t\*\*\*a\*\*\*e\*\*\*v\*\*\*

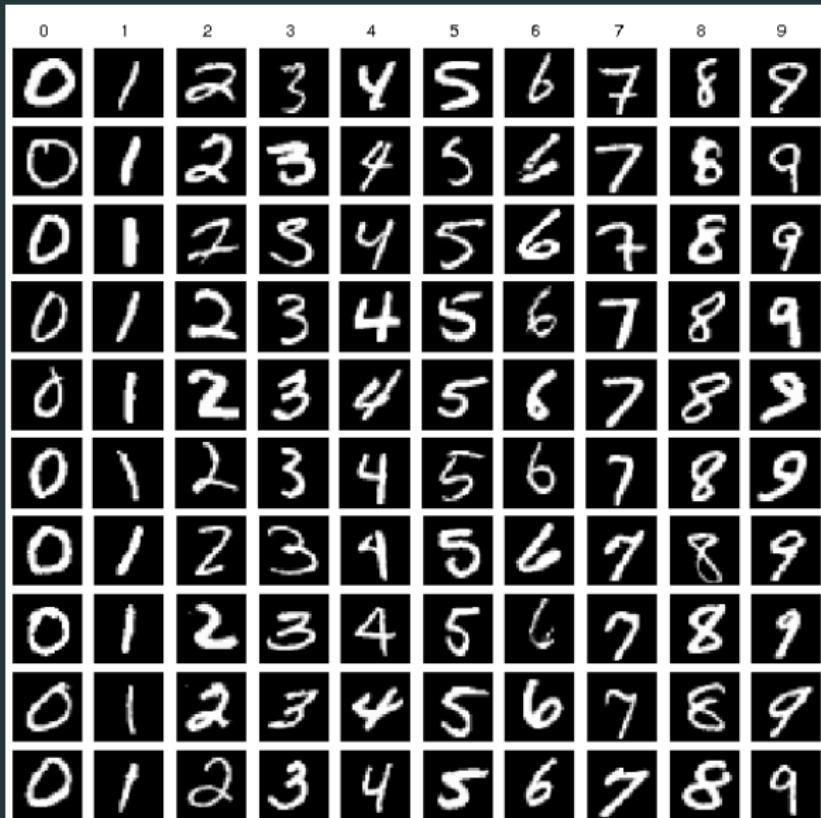
credit: David Mackay

Emma Woodhouse, handsome, clever, and rich, with a comfortable home and happy disposition, seemed to unite some of the best blessings of existence; and had lived nearly twenty one years in the world with very little to distress or vex her. She was the youngest of the two daughters of a most affectionate, indulgent father; and had, in consequence of her sister's marriage, been mistress of his house from a very early period. Her mother had died too long ago for her to have more than an indistinct remembrance of her caresses; and her place had been supplied by an excellent woman as governess, who had fallen little short of a mother in affection. Sixteen years had Miss Taylor been in Mr Woodhouse's family, less as a governess than a friend, very

## the noisy channel model in ML

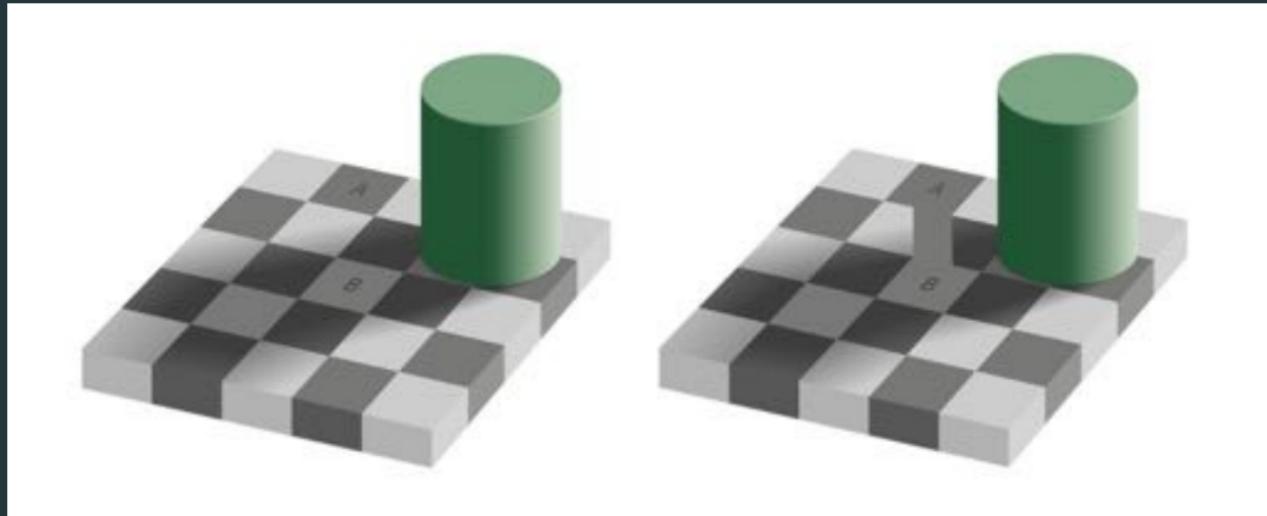


## noisy channels in ML



credit: MNIST dataset

## we are inherently bayesian



credit: quantamagazine.org, original image by Edward Adelson

“Tile A looks darker than tile B, though they are both the same shade (connecting the squares makes this clearer). The brain uses coloring of nearby tiles and location of the shadow to make inferences about the tile colors... lead to the perception that A and B are shaded differently.”

## what we hope to cover

- bayesian inference: unified paradigm for learning and decision-making
  - information theory: tool for designing and understanding data systems
- 
- basic probability review, and introducing information measures
  - the source and channel coding theorems

## what we hope to cover

- bayesian inference: unified paradigm for learning and decision-making
- information theory: tool for designing and understanding data systems

- basic probability review, and introducing information measures
- the source and channel coding theorems
- bayesian inference: priors, bayesian update, model selection
- generative models for discrete data: learning and inference

## what we hope to cover

- bayesian inference: unified paradigm for learning and decision-making
- information theory: tool for designing and understanding data systems

- basic probability review, and introducing information measures
- the source and channel coding theorems
- bayesian inference: priors, bayesian update, model selection
- generative models for discrete data: learning and inference
- bayesian graphical networks and markov random fields
- complex models: gaussian processes, artificial neural networks

## what we hope to cover

- bayesian inference: unified paradigm for learning and decision-making
- information theory: tool for designing and understanding data systems

- basic probability review, and introducing information measures
- the source and channel coding theorems
- bayesian inference: priors, bayesian update, model selection
- generative models for discrete data: learning and inference
- bayesian graphical networks and markov random fields
- complex models: gaussian processes, artificial neural networks
- approximate inference: MCMC and variational methods

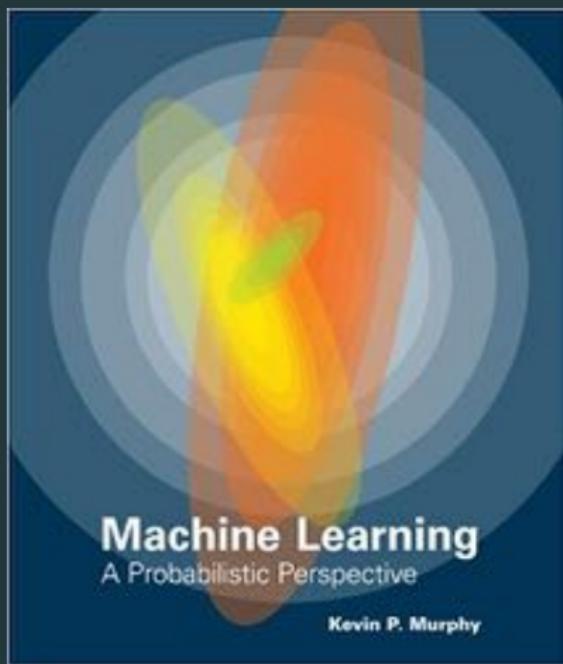
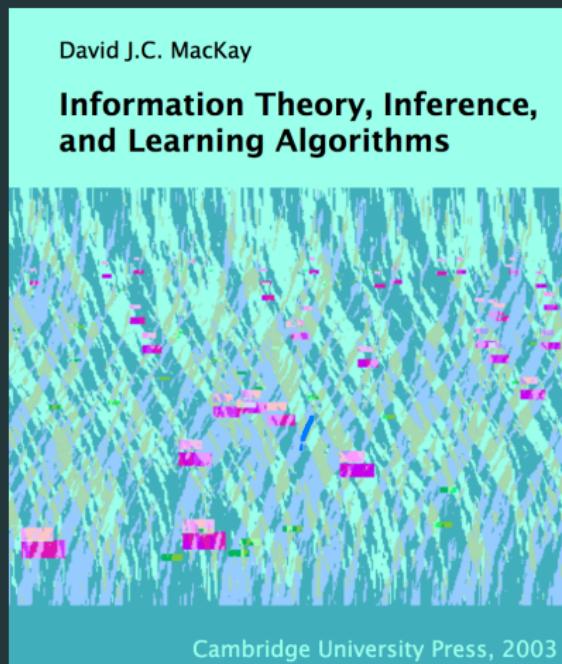
## what we hope to cover

- bayesian inference: unified paradigm for learning and decision-making
- information theory: tool for designing and understanding data systems

- basic probability review, and introducing information measures
- the source and channel coding theorems
- bayesian inference: priors, bayesian update, model selection
- generative models for discrete data: learning and inference
- bayesian graphical networks and markov random fields
- complex models: gaussian processes, artificial neural networks
- approximate inference: MCMC and variational methods
- model-based decision-making: bayesian optimization, causal inference, sequential decision-making and reinforcement learning

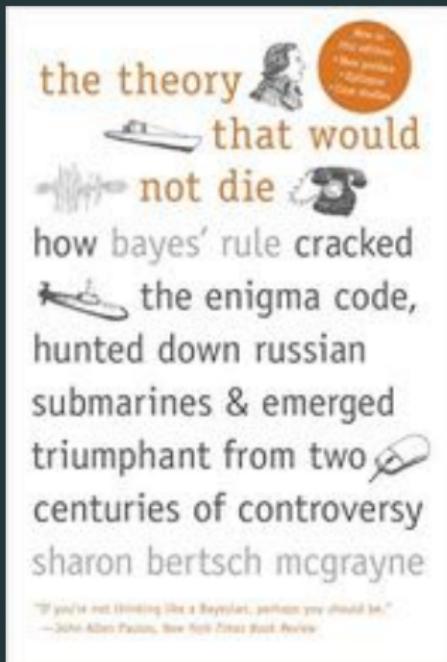
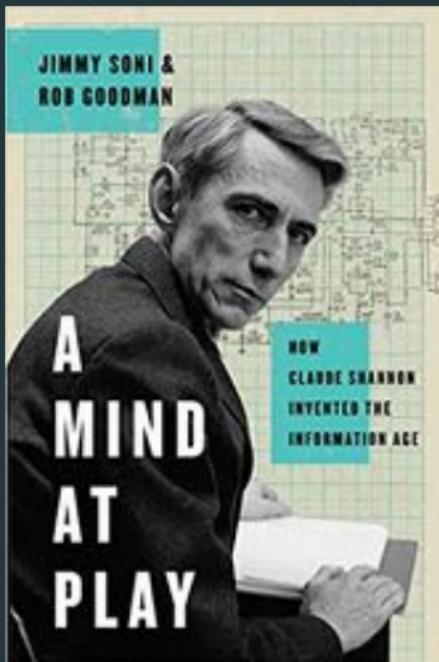
## aids in learning

the following books are excellent references for most topics in the course



# aids in getting excited about learning

the following help understand the larger context of what we will study



# is this course right for you?

- prerequisites:
  - linear algebra, calculus
  - probability: ideally at the level of ORIE 3500
  - programming: python/julia

# is this course right for you?

- prerequisites:
  - linear algebra, calculus
  - probability: ideally at the level of ORIE 3500
  - programming: python/julia
- caveat emptor:
  - may not be ideal as a first course in ML
  - we will focus on Bayesian methods, and ignore alternate 'frequentist' methods
  - will involve a fair bit of additional reading and programming, and some 'Bayesian philosophy'