



Bayesian statistics and Information Theory

Lecture 3: Information theory

... a.k.a. *how much is there to be learned in my data anyway?*

Florent Leclercq

www.florent-leclercq.eu

Imperial Centre for Inference and Cosmology

Imperial College London

May 28th, 2019

Outline: Lecture 3

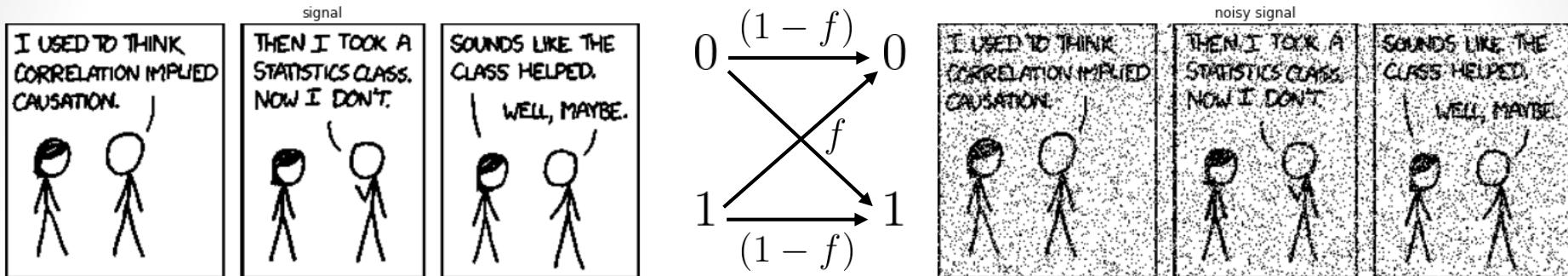
- Bayesian model comparison
 - Nested models and the Savage-Dickey density ratio
 - Bayesian model selection as a decision analysis
 - Bayesian model averaging
 - (Dangers of) model selection with insufficient summary statistics
- Information theory
 - The noisy binary symmetric channel
 - Low-density parity check codes
 - Measures of entropy and information
 - Information-theoretic experimental design
 - Supervised machine learning basics

Bayesian model comparison

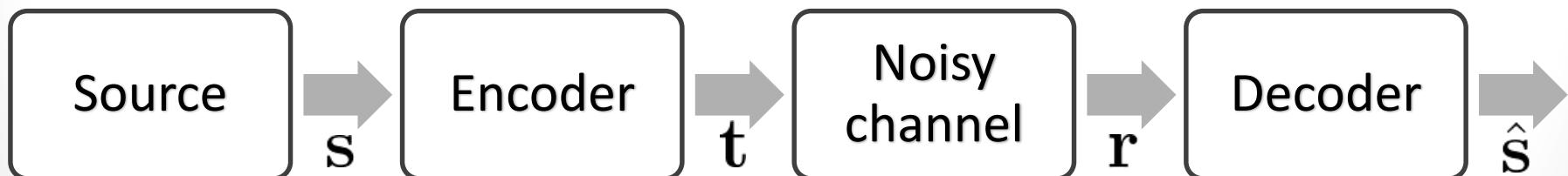
The noisy binary symmetric channel

Notebook 13: https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/IT_noisy_binary_channel.ipynb

The noisy binary symmetric channel



<https://xkcd.com/552/>



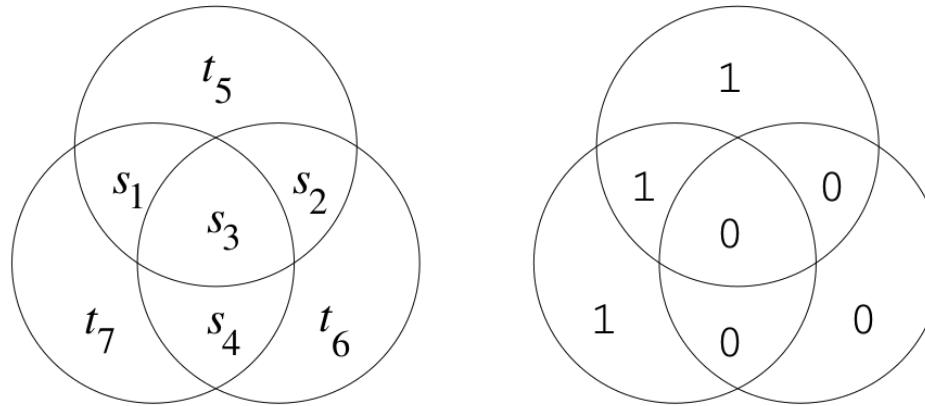
Rate of information transfer: $R = \frac{\#s}{\#t} = \frac{K}{N}$

The R₃ code

s	0	0	1	0	1	1	0
t	000	000	111	000	111	111	000
n	000	001	000	000	101	000	000
r	000	001	111	000	010	111	000
\hat{s}	0	0	1	0	0	1	0

The (7,4) Hamming code: Encoder

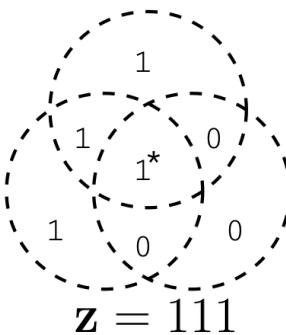
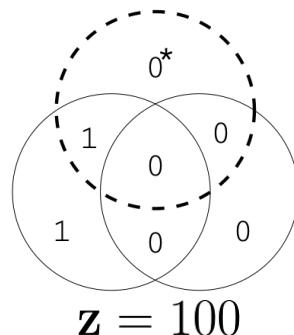
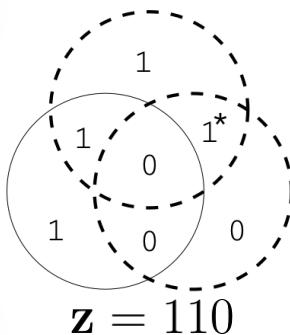
- Introducing the concept of **parity-checks**:



$$\mathbf{t} = \mathbf{Gs} \quad \mathbf{G} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

The (7,4) Hamming code: Decoder

- Introducing the concept of **syndrome**:
- Error correction:



$$\mathbf{H} = \left[\begin{array}{cccc|ccc} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 \end{array} \right] \quad \begin{matrix} \text{parity-checks} & & \text{identity} \end{matrix}$$

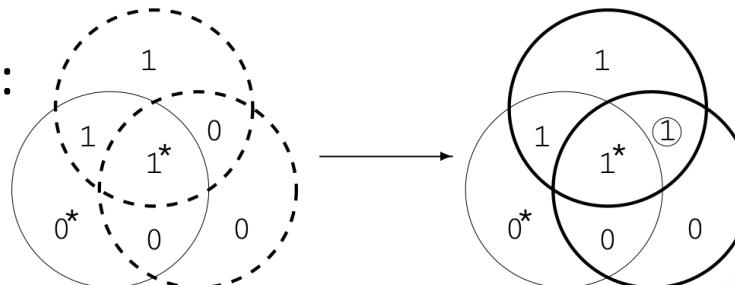
syndrome

$$z = \mathbf{Hr}$$

A Venn diagram showing three overlapping circles labeled r_1 , r_2 , r_3 , r_4 , r_5 , r_6 , and r_7 . The intersections of r_1 and r_2 , r_1 and r_3 , and r_2 and r_3 each contain a '1'. The intersection of all three circles contains a '0'. The regions outside all three circles but within the boundaries of the circles contain '1's. The region outside all seven circles contains a '0'. This corresponds to the syndrome $z = 110$.

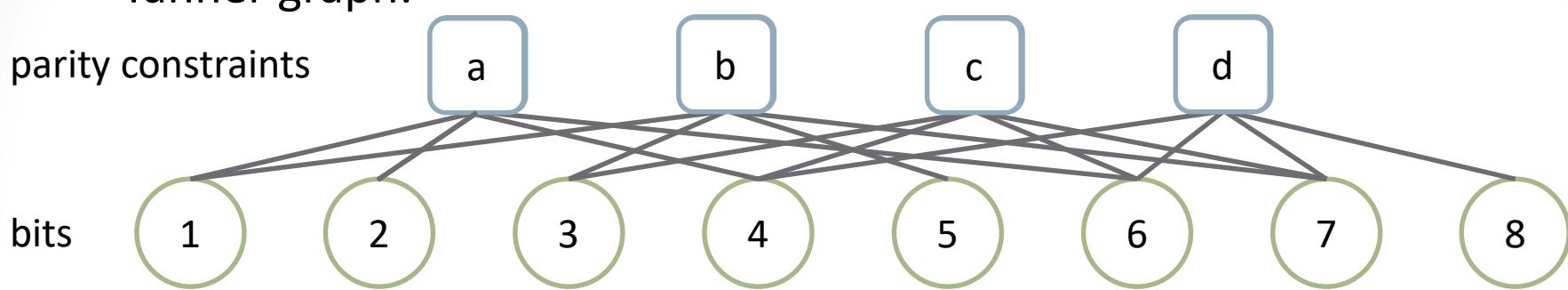
Syndrome z	000	001	010	011	100	101	110	111
Unflip this bit	None	r_7	r_6	r_4	r_5	r_1	r_2	r_3

- Unsuccessful error correction:



Low-density parity check (LDPC) codes

- Tanner graph:



- N bits, $M=(1-R)N$ parity-check constraints
 - 2^{RN} possible “words”: the *dictionary*
 - A *sparse* parity-check matrix
- Decoding LDPC codes: the general theory borrows from **statistical physics**: Ising spins in interaction and the BP mean field approximation (Bethe-Peierls – Belief Propagation).

Measures of entropy and information

Measures of entropy and information

- **Information content:** $I[X] \equiv - \sum_{x \in \mathcal{X}} \log_2 p(x)$
- **Entropy:**

$$H[X] \equiv \langle I[X] \rangle_{p(X)} = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x)$$

- **Conditional entropy:**

$$H[X|Y] \equiv \langle H[X|Y = y] \rangle_{p(Y)} = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log_2 \left(\frac{p(y)}{p(x, y)} \right)$$

- **Properties:**
 - chain rule: $H[X|Y] = H[X, Y] - H[Y]$
 - “Bayes’s theorem”: $H[X|Y] + H[Y] = H[Y|X] + H[X]$

Measures of entropy and information

- **Mutual information:** $I[X : Y] \equiv \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log_2 \left(\frac{p(x, y)}{p(x)p(y)} \right)$
- Properties:
 - $I[X : Y] \geq 0$
 - $I[X : X] = H[X]$
 - Consequence: $H[X|Y] \leq H[X]$
- **Relative entropy/Kullback-Leibler divergence/Information gain:**
 - Properties:
 - Gibbs's inequality: $D_{\text{KL}}[p||q] \geq 0$
 - Relation to mutual information:

$$D_{\text{KL}}[p||q] \equiv \sum_{x \in \mathcal{X}} p(x) \log_2 \left(\frac{p(x)}{q(x)} \right)$$

$$I[X : Y] = D_{\text{KL}}[p(x, y)||p(x)p(y)] = \langle D_{\text{KL}}[p(x|y)||p(x)] \rangle_{p(Y)}$$

Measures of entropy and information

- **Symmetrisation procedures:**

- Jeffreys symmetrisation: $C_s[A:B] = \frac{1}{2}C_a[A||B] + \frac{1}{2}C_a[B||A]$
- Jensen symmetrisation:

$$C_s[A:B] = \frac{1}{2}C_a[A||M] + \frac{1}{2}C_a[B||M] \quad \text{with} \quad M = \frac{A+B}{2}$$

- **Jeffreys divergence:** $D_J[p:q] = \frac{1}{2}D_{\text{KL}}[p||q] + \frac{1}{2}D_{\text{KL}}[q||p]$

- **Jensen-Shannon divergence:**

$$D_{\text{JS}}[p:q] = \frac{1}{2}D_{\text{KL}}[p||r] + \frac{1}{2}D_{\text{KL}}[q||r] \quad \text{with} \quad r \equiv \frac{p+q}{2}$$

- **Properties:**

- $D_{\text{JS}}[p:q] = H[r] - \frac{1}{2}H[p] - \frac{1}{2}H[q]$
- $0 \leq D_{\text{JS}}[p:q] \leq 1$
- $D_{\text{JS}}[p:q] = I[m:z] \quad \text{with} \quad m \equiv z p + (1-z) q \quad z = \begin{cases} 0 & (p = 1/2) \\ 1 & (p = 1/2) \end{cases}$

Information-theoretic experimental design

Information-theoretic experimental design

(expected) data experimental design

- **General problem:** maximize $U(\xi) = \langle U(d, \xi) \rangle_{p(d|\xi)} = \int p(d|\xi)U(d, \xi) \, dd$

- 1- Parameter inference utility functions:

- Maximal information gain:

$$U(d, \xi) \equiv D_{\text{KL}}[p(\theta|d, \xi) || p(\theta|\xi)] \quad U(\xi) = I[\theta : d | \xi]$$

- A-optimality:

$$U_A(d, \xi) \equiv \frac{1}{\text{tr}(\text{cov}(\theta|d, \xi)^{-1})}$$

- D-optimality:

$$U_D(d, \xi) \equiv \det(\text{cov}(\theta|d, \xi))$$

Information-theoretic experimental design

- **2- Model selection utility functions:**

- Maximal Bayes factor:

$$U(\xi) \equiv \mathcal{B}_{12}(\xi) = \frac{p(d|\xi, \mathcal{M}_1)}{p(d|\xi, \mathcal{M}_2)}$$

It is hard to predict Bayes factors...
but see [Trotta 2007, arXiv:astro-ph/0703063](#)

- Maximal mutual information between the model indicator and the data:

$$U(\xi) \equiv I[\mathcal{M}:d|\xi]$$

e.g. [Cavagnaro et al. 2010](#)

- Maximal Jensen-Shannon divergence between posterior predictive distributions:

$$U(\xi) \equiv D_{\text{JS}}[p_1:p_2] = I[\mathcal{M}:r|\xi] \quad \text{with} \quad r \equiv \frac{p_1 + p_2}{2}$$

[Vanlier et al. 2014, FL, Lavaux, Jasche & Wandelt 2016, arXiv:1606.06758](#)

- **3- Utilities for prediction of future observations:**

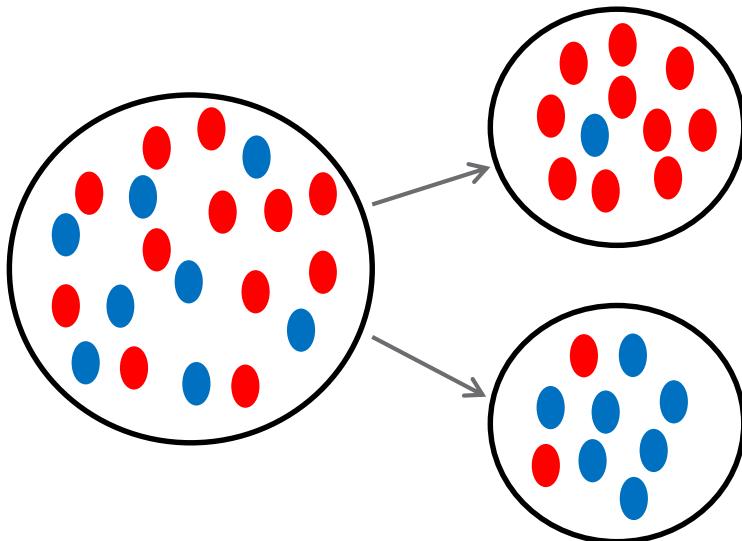
$$U(d, \xi) \equiv D_{\text{KL}}[p(t|d, \xi)||p(t|\xi)]$$

$$U(\xi) = I[t:d|\xi] = H[p(t|\xi)] - H[p(t|d, \xi)]$$

i.e. the “supervised machine learning” utility: $U(a) = H[T] - H[T|a]$

Supervised machine learning basics

- How to compute the information gain?



parent entropy:

$$H = -\frac{8}{20} \log_2 \left(\frac{8}{20} \right) - \frac{12}{20} \log_2 \left(\frac{12}{20} \right) = 0.9709$$

child1 entropy:

$$H = -\frac{10}{11} \log_2 \left(\frac{10}{11} \right) - \frac{1}{11} \log_2 \left(\frac{1}{11} \right) = 0.4395$$

child2 entropy:

$$H = -\frac{7}{9} \log_2 \left(\frac{7}{9} \right) - \frac{2}{9} \log_2 \left(\frac{2}{9} \right) = 0.7642$$

weighted average entropy of children:

$$\frac{11}{20} \times 0.4395 + \frac{9}{20} \times 0.7642 = 0.5856$$

information gain for this split: $0.9709 - 0.5856 = 0.3853$ Sh

Supervised machine learning: Titanic example

Notebook 14: https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/Machine_Learning_basics.ipynb