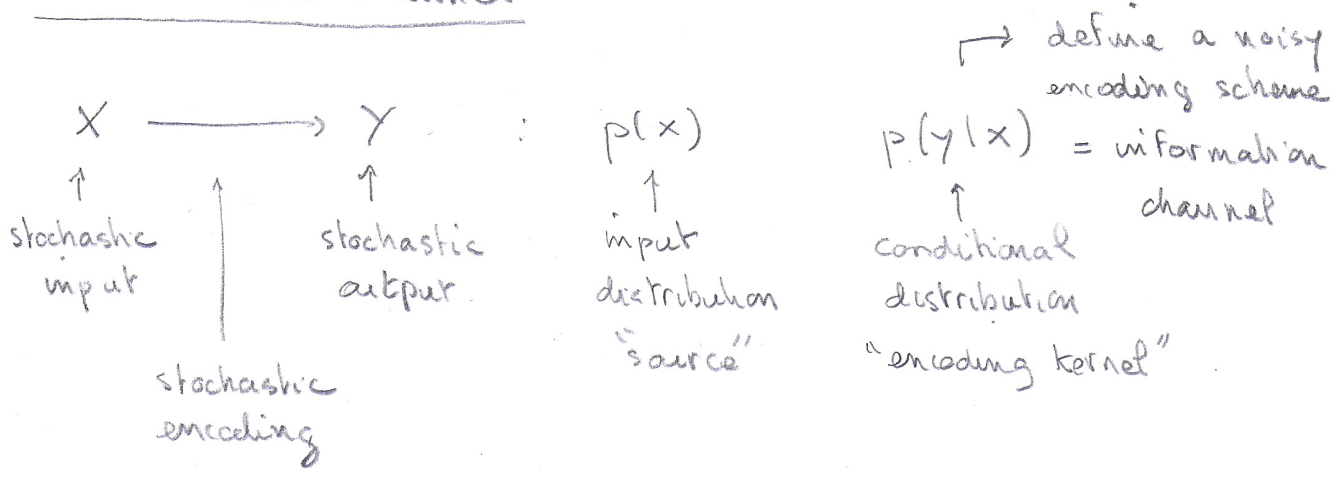


Information channel



In practice: * $p(x)$ can be thought of the probability of some stimulus feature that are relevant to behavior.

* $p(y|x)$ represents the noisy processing of stimulus information by neural network: for example, y is a vector of neuronal activity.

Efficient coding hypothesis: (debated) Neural system have evolved to optimize the flow/processing of information given biophysical constraints on the organism resource (range of firing rate, mean level of activity...)

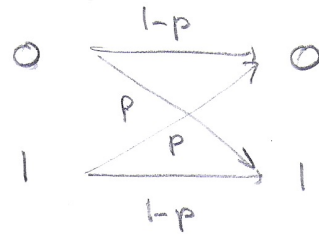
If relevant, the most important part would probably be to identify the nature of relevant constraints.

Problem : $\max I(X, Y)$ given some constraints on $p(x), p(y|x)$.

Channel capacity

Example: binary symmetric channel

p : error probability



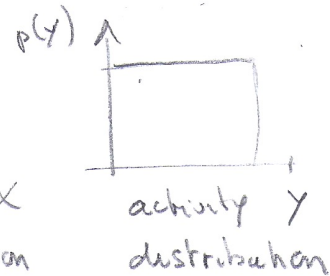
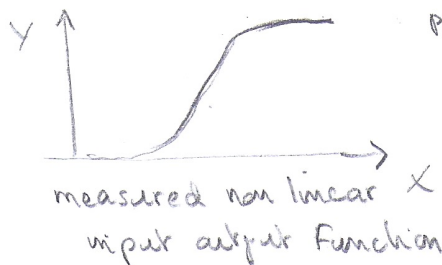
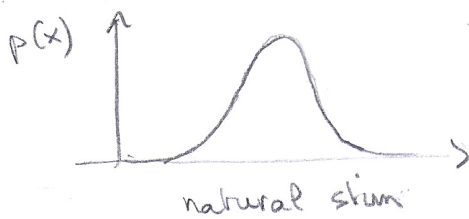
$$I(X, Y) = H(Y) - H(Y|X)$$

$$= H(Y) - \sum_x p(x) H(Y|X=x) \rightarrow H(p) \leftarrow \text{binary variable of parameter } p$$

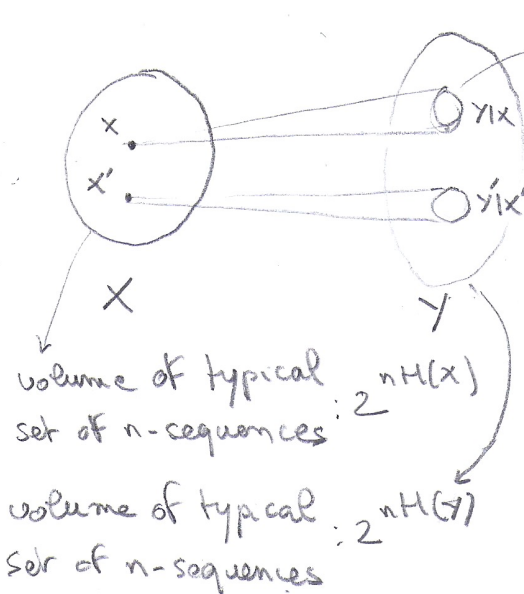
$$= H(Y) - H(p) \leq 1 - H(p) \leftarrow \text{maximum entropy of a binary variable is 1}$$

General principle: for encoding channel with $H(Y|X) = c$ the channel capacity corresponds to uniformly distributed output Y

in neuroscience: Laughlin's histogram equalization



Shannon's theorem:



average volume of typical set of n -sequences conditioned to a given output sequence: $2^{nH(Y|X)}$

$$\begin{aligned} \# \text{ of reliably transmittable sequences} \\ &= \frac{2^{nH(Y)}}{2^{nH(Y|X)}} = 2^{n(H(Y) - H(Y|X))} = 2^{nI(Y|X)} \end{aligned}$$

Capacity = optimal "soft" partitioning

Arimoto algorithm

(4)

↳ compute capacity and optimal input distribution

↳ idea: solve a bigger optimization problem with identical solution to find a monotonically converging method (MI increases each step)

$$I(x, y) = \sum_x \sum_y p(x) p(y|x) \log \frac{p(y|x)}{\sum_y p(x) p(y|x)} = I[p]$$

functional of p
↓
problem: non local term

Variational problem under constraint $\sum_x p(x) = 1, p(x) \geq 0 \forall x$

Localizing the expression to get exploitable KKT conditions?

$$I(x, y) = \sum_x \sum_y p(x) p(y|x) \log \frac{p(x|y)}{p(x)} \leftarrow \text{local}$$

but $p(x|y) = \frac{p(x) p(y|x)}{\sum_y p(x) p(y|x)}$ still non local \Rightarrow assume $p(x|y)$ unknown kernel denoted $q(x|y)$.

$$J[p, q] = \sum_x \sum_y p(x) p(y|x) \log \frac{q(x|y)}{p(x)} \leftarrow \text{concave in } p \text{ and } q!$$

Fact 1: $\max_p \max_q J[p, q] = C$

Proof: immediate if $\max_q J[p, q] = I[p]$

$$\text{KKT: } \frac{\delta}{\delta q(x|y)} \left[J[p, q] + \sum_y p(y) \left(\sum_x q(x|y) - 1 \right) \right] = 0$$

$$\Rightarrow \frac{p(y, x)}{q(x|y)} - p(y) = 0 \Rightarrow q(x|y) = \frac{p(y, x)}{p(y)} \underset{\substack{\uparrow \\ \text{normalization}}}{=} p(x|y)$$

✓

Fact 2: $\arg \max_p J[p, q] = p^*$, $p^*(x) = \frac{e^{\sum_y p(y|x) \log q(x|y)}}{\sum_x e^{\sum_y p(y|x) \log q(x|y)}}$

Proof: KKT: $\frac{\delta}{\delta p[x]} [J[p, q] - \mu (\sum p(x) - 1)] = 0$

$$\sum_y p(y|x) \left[\log \left(\frac{q(x|y)}{p(x)} \right) - 1 \right] - \mu = 0$$

$$\log p(x) = \sum_y p(y|x) \log q(x|y) - \mu + 1 \quad \checkmark$$

Algorithm: iterates two optimization steps:

geometric interpretation? $\left\{ \begin{array}{l} q^{(n)} \leftarrow q^{(n+1)} / q^{(n+1)}(x) \propto p^{(n)}(x) p(y|x) \\ p^{(n)} \leftarrow p^{(n+1)} / p^{(n+1)}(x) \propto \exp \left[\sum_y p(y|x) \log q^{(n+1)}(x|y) \right] \end{array} \right.$

it can be shown that iterating these steps leads to a strictly increasing sequence of MI values converging to the capacity C .

At capacity $p(x) \propto \exp \left[\sum_y p(y|x) \log \left(\frac{p(x) p(y|x)}{\sum_y p(x) p(y|x)} \right) \right]$

$$p(x) \propto p(x) D_{KL}(p(y|x) || p(y)) \rightarrow = \text{constant}$$

$p^*(x)$ is such that $p^*(y)$ is equally distinct from each the conditional distribution $p(y|x)$ as measured by Kullback-Leibler divergence

Minimax interpretation

↳ setting from game theory

P_θ
" "

- 1 - nature picks θ from a statistical model $p(x|\theta), \theta \in \Theta$
- 2 - statistician guesses P_θ by choosing some distribution q

Loss function is measured by $D_{KL}(p(x|\theta) || q(x))$

Game: nature wants to find a prior $p(\theta)$ that maximize the statistician loss, while the statistician want a decision strategy that minimizes his loss.

Statistician: $\arg \min_q \int D_{KL}(p(x|\theta) | q(x)) p(\theta) d\theta = q^*$

$q^*(x) = \int p(x|\theta) p(\theta) d\theta \rightarrow$ Bayes strategy

Nature: $\max_P \left(\min_q \left(\int D_{KL}(p(x|\theta) | q(x)) p(\theta) d\theta \right) \right) = C$

maximin strategy \uparrow $I(x, \theta)$ \uparrow capacity

What about minimax strategies?

↳ It turns out that the min and max commute (Hausler 1995)

$C = \min_q \left(\max_P \left(\int D_{KL}(p(x|\theta) || q(x)) p(\theta) d\theta \right) \right) -$
 \uparrow
 linear in $p(\theta)$!

Rate-distortion

↳ achieving high information rate is "costly" in terms of coding (redundant codes) to average noise

↳ idea: looking for least informative encoding given some constraints on the faithfulness of the encoding, as measure by a distortion metric, e.g. $d(x, y) \geq 0$ $d(x, x) = 0$

Rate-distortion function:

non local

$$R(\delta) = \min_{p(y|x)} I(x, y) = \min_{p(y|x)} \sum_x \sum_y p(x) p(y|x) \log \frac{p(y|x)}{\sum_x p(y|x) p(x)}$$

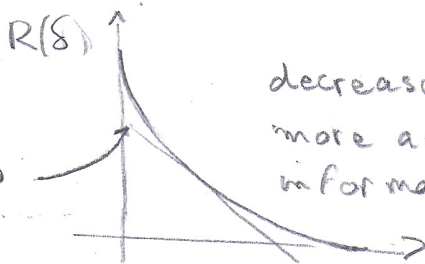
under average distortion constraint $\sum_x \sum_y p(x) p(y|x) d(x, y) \leq \delta$

Lagrange multiplier s associated to constraint δ

$$R(\delta) = \min_{p(y|x)} \left[I(x, y) - s \sum_x \sum_y p(x) p(y|x) d(x, y) \right]$$

Intuitively

$s \leq 0 =$ slope of R at δ



decreasing function of δ : more accuracy ($\delta \downarrow$) requires more information to be transmitted

Same localization idea as for the capacity introducing

$$F[p|q] = \sum_x p(x) p(y|x) \log \frac{p(y|x)}{q(y)} - s \sum_x \sum_y p(x) p(y|x) d(x, y)$$

$$R[\delta] = \min_p \min_q F[p|q] \Rightarrow q^{(n)}(y) \leftarrow q^{(n+1)}(y) = \sum_x p(x) p^{(n)}(y|x)$$

Blahut algorithm

$$p^{(n+1)}(y|x) \leftarrow p^{(n)}(y|x) \propto q^{(n+1)}(y) e^{s d(x, y)}$$