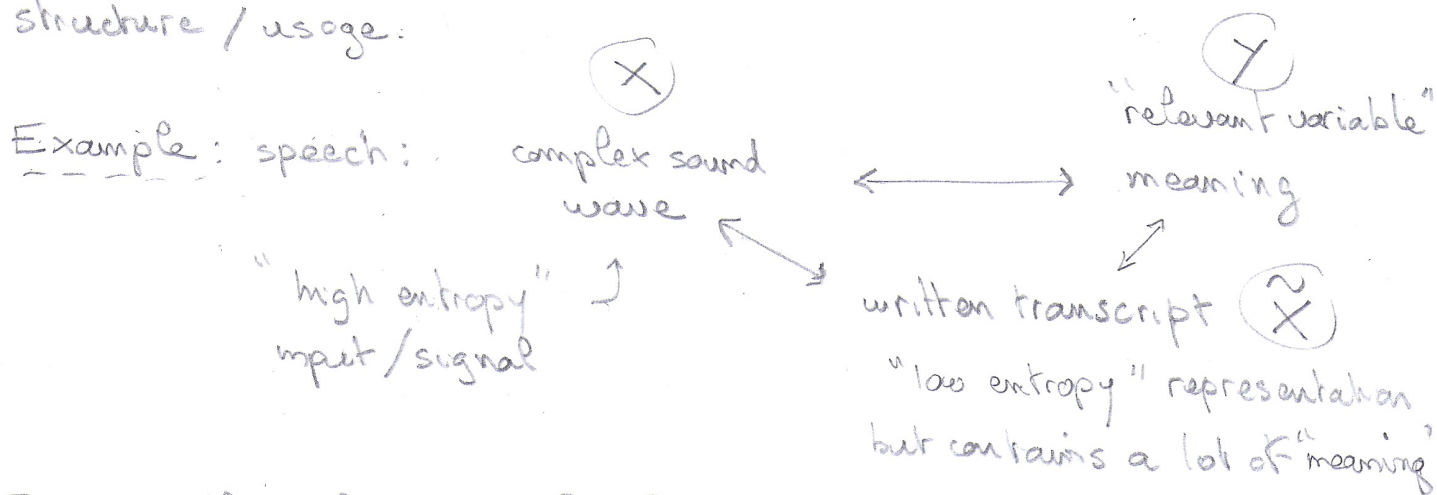


Information relevance

①

The original formulation of information theory by Shannon purposely ignored the problem of information relevance, i.e. judging the value - or meaning - of information to the recipient. Only concern: transmission efficiency

Problem of information relevance is a priori not well-posed in the absence of additional knowledge about information structure / usage:



Idea: the relevance of information conveyed about X in the representation \tilde{X} can be measured via the introduction of an additional relevance variable Y

Tradeoff: there is a natural tradeoff between the representation size and the expected distortion of the meaningful content. \longleftrightarrow akin to rate-distortion theory

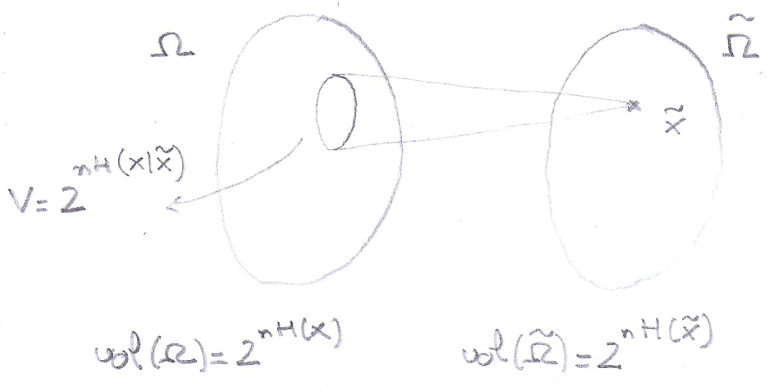
Choice of relevant variable \longleftrightarrow Feature selection

Rate distortion theory

$X \rightarrow \tilde{X}$: encoding channel

X : input / signal / stimulus
 \tilde{X} : code / pattern / representation

↑
stochastic mapping $p(\tilde{X}|X)$ = soft partitioning / quantization



V = average volume of input space that get mapped to the same code.

distinguishable states
 $= 2^{nH(X)} / 2^{nH(\tilde{X}|X)} = 2^{nI(X, \tilde{X})}$

Quality of the quantization is measured by the rate of bits needed to specify input without confusion, which is bounded below by $I(X, \tilde{X})$.

Trade off between rate and quality of the reconstruction as measured via a distortion function:

$d: X \times \tilde{X} \rightarrow \mathbb{R}^+$: d small when x can be faithfully reconstructed from \tilde{x} .

$\mathbb{E}_{p(x, \tilde{x})} [d(x, \tilde{x})]$: expected distortion given a code

↳ the larger the rate, the smaller the achievable distortion

Convex optimization setting:

$R(D) = \min_{p(\tilde{X}|X) \text{ s.t. } \mathbb{E}[d] \leq D} I(X, \tilde{X})$ ← minimal achievable rate under constraint of expected distortion

Variational problem and iterative algorithm

$$F[p(\tilde{x}|x)] = I(x, \tilde{x}) + \beta \mathbb{E}[d(x, \tilde{x})] + \lambda(x) \sum p(\tilde{x}|x)$$

\uparrow Lagrangian \uparrow Lagrangian multiplier \uparrow normalization

$$F[p(\tilde{x}|x)] = \sum_{x, x'} p(x) p(\tilde{x}|x) \log p(\tilde{x}|x) - \sum_{\tilde{x}} p(\tilde{x}) \log p(\tilde{x}) + \beta \sum_x p(x) \sum_{\tilde{x}} p(\tilde{x}|x) d(x, \tilde{x}) + \lambda(x) \sum_{\tilde{x}} p(\tilde{x}|x)$$

\uparrow non-local term

$$\frac{\delta F}{\delta p(\tilde{x}|x)} = p(x) (1 + \log p(\tilde{x}|x)) - \sum_{\tilde{x}'} \frac{\delta p(\tilde{x}')}{\delta p(\tilde{x}|x)} (1 + \log p(\tilde{x}')) + \beta p(x) d(x, \tilde{x}) + \lambda(x)$$

$\hookrightarrow = \delta_{\tilde{x}\tilde{x}'} p(x')$

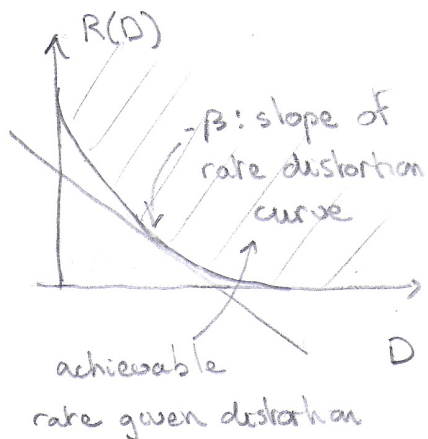
assumed non-zero \leftarrow

$$= p(x) \left(1 + \log p(\tilde{x}|x) - 1 - \log p(\tilde{x}) + \beta d(x, \tilde{x}) + \frac{\lambda(\tilde{x})}{p(\tilde{x})} \right)$$

$$\frac{\delta F}{\delta p(\tilde{x}|x)} = 0 \Rightarrow p(\tilde{x}|x) = \frac{p(\tilde{x})}{Z(x, \beta)} e^{-\beta d(x, \tilde{x})}$$

$Z(x, \beta) \leftarrow$ normalization constant "Free energy"

Moreover, at optimum: $\delta F = \delta I + \beta \delta \mathbb{E}[d] = 0 \Rightarrow \beta = - \frac{\delta R}{\delta D}$



self consistent equations

(*) $p(\tilde{x}) = \sum_x p(\tilde{x}|x) p(x)$

(**) $p(\tilde{x}|x) = \frac{p(\tilde{x}) e^{-\beta d(x, \tilde{x})}}{Z(x, \beta)}$

Iterative Blahut algorithm: $(*) \overset{\text{step } n+1}{p_{n+1}(\tilde{x})} = \sum_x \overset{\text{step } n}{p_n(\tilde{x}|x)} p(x)$

$(**) p_{n+1}(\tilde{x}|x) = \frac{p_{n+1}(\tilde{x}) e^{-\beta d(x, \tilde{x})}}{Z(x, \beta)}$

Justification:

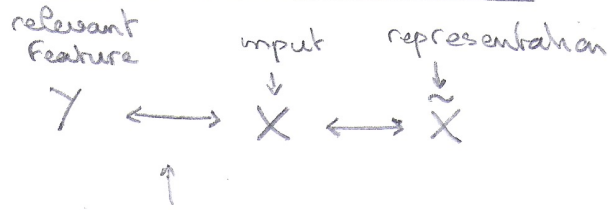
$$R(D) = \min_{p(\tilde{x}|x) \text{ s.t. } \mathbb{E}[d] \leq D} I(x, \gamma) = \min_{p(\tilde{x}|x)} \min_{p(\tilde{x})} F(p(\tilde{x}|x), p(\tilde{x})) \text{ s.t. } \mathbb{E}[d] \leq D$$

$$F(p(\tilde{x}|x), p(\tilde{x})) = \sum_{\tilde{x}, x} p(x) p(\tilde{x}|x) \log \frac{p(\tilde{x}|x)}{p(\tilde{x})}$$

↑
 minimization of F can be performed independently on the convex sets of probabilities for $p(\tilde{x})$ and $p(\tilde{x}|x), x \in X$

Information bottleneck

(5)



$\tilde{X}|X$ and $X|X$ are independent

↳ Markov information channel

⚠ double arrow notation to emphasize lack of directionality

$$Y \rightarrow X \rightarrow \tilde{X} \approx Y \leftarrow X \leftarrow \tilde{X}$$

- * Relevant representations are those for which $I(\tilde{X}, Y)$ is high. Remember that $I(\tilde{X}, Y) \leq I(X, Y)$ by the data processing inequality.
- * There is a trade-off between preserving relevant information and compressing the representation. The goal is to pass as much information about Y via X , but through the "bottleneck" formed by compact representations \tilde{X} .
- * Variational formulation: $L[p(\tilde{x}|x)] = I(X, \tilde{X}) - \beta I(X, Y)$

↑
Lagrangian multiplier

$\beta = 0$: most sketchy representation

$\beta \rightarrow \infty$: arbitrary detailed quantization

↳ no explicit distortion function!

↳ nonlinear convex optimization problem!

Self consistent equations

(6)

$$L[p(\tilde{x}|x)] = \sum_{x, \tilde{x}} p(x) p(\tilde{x}|x) \log p(\tilde{x}|x) - \sum_{\tilde{x}} p(\tilde{x}) \log p(\tilde{x}) \leftarrow \text{non-local}$$

$$- \beta \sum_{y, \tilde{x}} p(y) p(\tilde{x}|y) \log p(\tilde{x}|y) + \beta \sum_y p(\tilde{x}) \log p(\tilde{x}) \leftarrow$$

$$- \sum_x \lambda(x) \sum_{\tilde{x}} p(\tilde{x}|x) \leftarrow \text{normalization constraints}$$

$$\frac{\delta L}{\delta p(\tilde{x}|x)} = p(x) (1 + \log p(\tilde{x}|x)) - \frac{\delta p(\tilde{x})}{\delta p(\tilde{x}|x)} (1 + \log p(\tilde{x}))$$

$$- \beta \sum_y p(y) \frac{\delta p(\tilde{x}|y)}{\delta p(\tilde{x}|x)} (1 + \log p(\tilde{x}|y)) + \beta \frac{\delta p(\tilde{x})}{\delta p(\tilde{x}|x)} (1 + \log p(\tilde{x}))$$

$$- \lambda(x)$$

Markov chain $Y \leftrightarrow X \leftrightarrow \tilde{X} :$

$$p(\tilde{x}) = \sum_x p(\tilde{x}|x) p(x)$$

$$p(\tilde{x}|y) = \sum_x p(\tilde{x}|x) p(x|y)$$

$$\frac{\delta p(\tilde{x})}{\delta p(\tilde{x}|x)} = p(x) \quad \text{and} \quad \frac{\delta p(\tilde{x}|y)}{\delta p(\tilde{x}|x)} = p(x|y)$$

$$\frac{\delta L}{\delta p(\tilde{x}|x)} = p(x) \left\{ \begin{aligned} & 1 + \log p(\tilde{x}|x) - 1 - \log p(\tilde{x}) - \beta \sum_y p(y|x) (1 + \log p(\tilde{x}|y)) \\ & + \beta (1 + \log p(\tilde{x})) - \frac{\lambda(x)}{p(x)} \end{aligned} \right\}$$

assumed $\neq 0$ \uparrow

$$= p(x) \left\{ \log \frac{p(\tilde{x}|x)}{p(\tilde{x})} - \beta \sum_y p(y|x) \log \frac{p(y|\tilde{x})}{p(y)} - \frac{\lambda(x)}{p(x)} \right\}$$

$$= p(x) \left\{ \log \frac{p(\tilde{x}|x)}{p(\tilde{x})} + \beta \sum_y p(y|x) \log \frac{p(y|x)}{p(y|\tilde{x})} \right\}$$

$D_{KL}(p(y|x) || p(y|\tilde{x}))$
"def"
 $d(x || \tilde{x})$

new multiplier: $\tilde{\lambda}(x) \leftarrow - \left(\frac{\lambda(x)}{p(x)} + \beta \sum_y p(y|x) \log \frac{p(y|x)}{p(y)} \right)$

$$\frac{\delta L}{\delta p(\tilde{x}|x)} = 0 \Rightarrow p(\tilde{x}|x) = \frac{p(\tilde{x})}{Z(x, \beta)} e^{-\beta d(x, \tilde{x})}$$

↑ normalization
↑ constant
↑ distortion function

The Kullback-Leibler divergence $d(x, \tilde{x}) = D_{KL}(p(y|x) \| p(y|\tilde{x}))$ emerges as the relevant distortion function.

At optimum: self consistent equations

$$(*) \quad p(\tilde{x}) = \sum_x p(\tilde{x}|x) p(x)$$

$$(**) \quad p(y|\tilde{x}) = \sum_x p(y|x) p(x|\tilde{x})$$

$$(***) \quad p(\tilde{x}|x) = \frac{p(\tilde{x})}{Z(x, \beta)} e^{-\beta d(x, \tilde{x})}$$

Markov chain
 $Y \leftrightarrow X \leftrightarrow \tilde{X}$

Iterative algorithm

$$1. p_n(\tilde{X}|X) = \frac{p_n(\tilde{x})}{Z_n(x, \beta)} e^{-\beta D_{KL}(p(y|x) || p_n(y|x))} \leftarrow (***)$$

$$2. p_{n+1}(\tilde{X}) = \sum_x p_n(\tilde{X}|x) p(x) \leftarrow (*)$$

$$2'. p_{n+1}(Y|\tilde{X}) = \sum_x p(y|x) p_n(x|\tilde{X}) \leftarrow (**)$$

Interpretation

$$\min_{p(\tilde{x}|x) \text{ s.t. } I(\tilde{X}, Y) \geq D} I(X, \tilde{X}) = \min_{p(y|\tilde{x})} \min_{p(x)} \min_{p(\tilde{x}|x)} F[p(\tilde{x}|x), p(x), p(y|\tilde{x})]$$

s.t. $I(\tilde{X}, Y) \geq D$

$$F[p(\tilde{x}|x), p(x), p(y|\tilde{x})] = I(X, \tilde{X}) + \beta \mathbb{E}_{p(x, \tilde{x})} [D_{KL}(p(y|x) || p(y|\tilde{x}))]$$

$+ K[X, Y] \leftarrow \text{constant}$

regular rate-distortion
Function minimized when
(*) and (***) hold at fixed
 $p(y|\tilde{x})$

Functional of $p(y|\tilde{x})$
minimized when (**)
holds at fixed $p(x, \tilde{x})$