

The Information Bottleneck Method was proposed by Tishby, Pereira, and Bialek (Allerton, 1999). Lecture based on Goldkila & Polyanskiy (2020) and Zaidi et al (2020).

An information-theoretic framework for learning.

Considers extracting information about a target signal  $Y$  through a correlated observable  $X$ .

The extracted information is quantified by a variable  $T = T(X)$  which is a possibly randomized function of  $X$ , thereby forming Markov chain

$$Y \rightarrow X \rightarrow T.$$

Objective is to find a  $T$  that minimizes the mutual information  $I(X; T)$  while keeping  $I(Y; T)$  above a certain threshold.

Threshold determines how informative the representation  $T$  is about  $Y$ .

$$\begin{aligned} \min I(X; T) \\ \text{s.t. } I(Y; T) \geq \alpha. \end{aligned}$$

where minimization over all randomized mappings of  $X$  to  $T$ .

This was originally motivated axiomatically rather than with

a coding theorem, or a precise operational problem: want to keep as much relevant/informative info as possible.  
 $\rightarrow$  leads to a clustering algorithm.

IB framework concerned with finding  $P_{T|X}$  that extracts information about  $Y$ , i.e. high  $I(Y;T)$ , while compressing  $X$ , quantified as keeping  $I(X;T)$  small.

Data processing inequality in information theory for Markov chain  $Y \rightarrow X \rightarrow T$  implies  $I(Y;T) \leq I(X;T)$ ,

so compressed representation  $T$  can't convey more information than original signal.

so, IB basically trying to find approximate version of minimal sufficient statistic

We pass the information  $X$  contains about  $Y$  through a bottleneck via the representation  $T$ .

$\Rightarrow$  claim is that this objective promotes minimality, sufficiency, and disentanglement of representations.

A slightly different form of the same ~~formulation~~ problem

$$\begin{aligned} \min \quad & H(Y|T) \\ P_{T|X} : & H(X|T) \geq \alpha \end{aligned}$$

was developed by Witsenhausen and Wyner (1975), in context of common information.

The original problem statement not convex, but can be made into a convex rate-distortion problem. Commonly solved by introducing Lagrange multiplier  $\beta$  and considering the functional

$$J_\beta(P_{T|X}) = I(X;T) - \beta I(T;Y).$$

Here  $\beta$  controls amount of compression

→ small  $\beta$  implies more compression

→ large  $\beta$  leads to more informativeness/relevance.

The Lagrangian formulation then leads to a variational characterization of stationary point of  $J_\beta(P_{T|X})$

A stationary point  $P_{T|X}^{(\beta)}$  must satisfy

$$\textcircled{1} \quad P_T^{(\beta)}(t) = \int_{\mathcal{X}} P_{T|X}^{(\beta)}(t|x) dP_X(x) \quad (\text{this is Radon-Nikodym derivative})$$

$$\textcircled{2} \quad P_{Y|T}^{(\beta)}(y|t) = \frac{1}{P_T^{(\beta)}(t)} \int_{\mathcal{X}} P_{Y|X}(y|x) P_{T|X}^{(\beta)}(t|x) dP_X(x)$$

$$\textcircled{3} \quad P_{T|X}^{(\beta)}(t|x) = \frac{P_T^{(\beta)}(t)}{Z_\beta(x)} e^{-\beta D_{KL}(P_{Y|X}(\cdot|x) || P_{Y|T}^{(\beta)}(\cdot|t))}$$

where  $Z_\beta(x)$  is normalization constant (partition function).

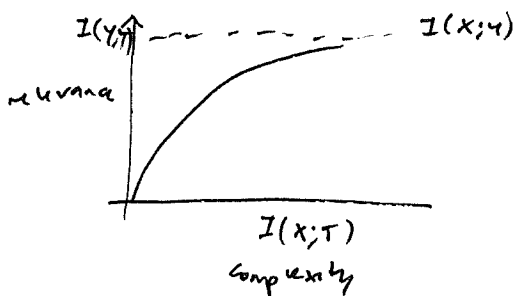
If  $X, Y, T$  take values in finite sets, and  $P_{X,Y}$  known, then alternating

iterations of  $\left( \begin{matrix} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \end{matrix} \right)$  locally converge to a solution for any initial  $P_{T|X}$ .

Reminiscent of Blahut-Arimo algorithm for rate distortion.

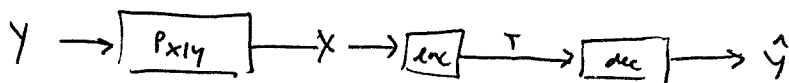
The IB curve is obtained by plotting

$(I_B(x;T), I_B(y;T))$  for an optimal  $P_{T|X}^{(A)}$  for each  $\beta \in [0, \infty)$ .



Come back to operational question of what problem is the information bottleneck the solution to? Looks a lot like a rate distortion problem with a particular information-based distortion function.

Consider the remote source coding



with a distortion measure that is the log-loss fidelity criterion.

The decoder generates a soft estimate  $\hat{y}^n$  of  $y^n$  in the form of a probability distribution over  $\mathcal{Y}^n$ , i.e.  $\hat{y} = \hat{P}_{Y^n|M}(\cdot)$ .

The incurred discrepancy between  $y^n$  and the estimate  $\hat{y}^n$  under log-loss for the observation  $x^n$  is then the per-letter log-loss distortion:

$\ell_{\log}(y, \hat{y}) = \log \frac{1}{\hat{g}(y)}$  for  $y \in \mathcal{Y}$  and  $\hat{g} \in \mathcal{P}(\mathcal{Y})$  designates  
 a probability distribution on  $\mathcal{Y}$  and  $\hat{g}(y)$  is value of that  
 distribution evaluated at outcome  $y \in \mathcal{Y}$ .

---

So want  $E[\ell_{\log}^{(n)}(Y^n, \hat{Y}^n)] \leq D$

where incurred distribution between sequences  $Y^n$  and  $\hat{Y}^n$  is measured as

$$\ell_{\log}^{(n)}(Y^n, \hat{Y}^n) = \frac{1}{n} \sum_{i=1}^n \ell_{\log}(y_i, \hat{y}_i).$$

It is known, rate-distortion region is  $(R, D)$  that satisfy

$$R \geq \underline{I}(U; X)$$

$$D \geq H(Y|U)$$

where ~~choice~~ <sup>allowed</sup> is over all auxiliary random variables satisfying

$$U \rightarrow X \rightarrow Y, \text{ and } |U| \leq |X| + 1.$$

Using substitution  $\Delta = H(Y) - D$ , get pairs  $(R, H(Y) - D)$  that

satisfy  $R \geq \underline{I}(U; X)$

$$\Delta \leq \underline{I}(U; Y)$$

for  $U$  that satisfies  $U \rightarrow X \rightarrow Y$  with  $|U| \leq |X| + 1$ .

Also connected to WAK problem.

---

Sources  $X$  and  $Y$  encoded separately at rates  $R_x$  and  $R_y$   
 decoder recovers  $Y$  losslessly.

For given  $R_x = R$ , minimum rate  $R_y$  to losslessly recover  $Y$  is

$$R_y^*(R) = \min_{P_{U|X}: I(U;X) \leq R} H(Y|U)$$

so we get

$$\max_{P_{U|X}: I(U;X) \leq R} I(U;Y) = H(Y) - R_y^*(R).$$

so solving IB problem equivalent to solving WAK problem.