

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2020

Noisy Channel RDAs

Noisy-Channel RDAs

In the image compression case study, training was based on a differentiable loss

$$\Phi^* = \operatorname{argmin}_{\Phi} \left(E_y - \ln p_{\Phi}(z_{\Phi}(y)) \right) + \lambda E_{y,\epsilon} \operatorname{Dist}(y, y_{\Phi}(z_{\Phi}(y) + \epsilon))$$

In a rate-distortion auto-encoder we will replace the conceptually dubious differential entropy rate term with a conceptually legitimate mutual information (channel capacity) rate term.

Mutual Information as a Channel Rate

$$\Phi^* = \operatorname{argmin}_{\Phi} \left(E_y - \ln p_{\Phi}(z_{\Phi}(y)) \right) + \lambda E_{y,\epsilon} \operatorname{Dist}(y, y_{\Phi}(z_{\Phi}(y) + \epsilon))$$

is replaced by

$$\tilde{z} = z_{\Phi}(y) + f_{\Phi}(y, \epsilon) \quad \epsilon \text{ is fixed (parameter independent) noise}$$

$$\Phi^* = \operatorname{argmin}_{\Phi} I_{\Phi}(y, \tilde{z}) + \lambda E_{y,\epsilon} \operatorname{Dist}(y, y_{\Phi}(\tilde{z}))$$

By the channel capacity theorem $I(y, \tilde{z})$ is the **rate** of information transfer from y to \tilde{z} .

Mutual Information as a Channel Rate

$\tilde{z} = z_{\Phi}(y) + f_{\Phi}(y, \epsilon)$ ϵ is fixed (parameter independent) noise

$$\Phi^* = \operatorname{argmin}_{\Phi} I_{\Phi}(y, \tilde{z}) + \lambda E_{y, \epsilon} \operatorname{Dist}(y, y_{\Phi}(\tilde{z}))$$

Taking the distribution on ϵ to be parameter independent is called the “reparameterization trick” and allows SGD.

$$\begin{aligned} & \nabla_{\Phi} E_{y, \epsilon} \operatorname{Dist}(y, y_{\Phi}(z_{\Phi}(y) + f_{\Phi}(y, \epsilon))) \\ &= E_{y, \epsilon} \nabla_{\Phi} \operatorname{Dist}(y, y_{\Phi}(z_{\Phi}(y) + f_{\Phi}(y, \epsilon))) \end{aligned}$$

Mutual Information as a Channel Rate

$\tilde{z} = z_{\Phi}(y) + \textcolor{red}{f}_{\Phi}(y, \epsilon)$ ϵ is fixed (parameter independent) noise

$$\Phi^* = \operatorname{argmin}_{\Phi} I_{\Phi}(y, \tilde{z}) + \lambda E_{y, \epsilon} \operatorname{Dist}(y, y_{\Phi}(\tilde{z}))$$

Typically $f_{\Phi}(y, \epsilon)$ is simple, such as $\sigma_{\Phi}(y) \odot \epsilon$, so that $\textcolor{red}{p}_{\Phi}(\tilde{z}|y)$ is easily computed.

Mutual Information Replaces Cross Entropy

$$I_{\Phi}(y, \tilde{z}) = E_{y, \epsilon} \ln \frac{\text{pop}(y) p_{\Phi}(\tilde{z}|y)}{\text{pop}(y) p_{\text{pop}, \Phi}(\tilde{z})}$$

$$= E_{y, \epsilon} \ln \frac{p_{\Phi}(\tilde{z}|y)}{p_{\text{pop}, \Phi}(\tilde{z})}$$

where $p_{\text{pop}, \Phi}(\tilde{z}) = E_{y \sim \text{pop}} p_{\Phi}(\tilde{z}|y)$

A Variational Bound

$$p_{\text{pop},\Phi}(\tilde{z}) = E_{y \sim \text{pop}} p_{\Phi}(\tilde{z}|y)$$

We cannot compute $p_{\text{pop},\Phi}(\tilde{z})$.

Instead we will use a variational bound involving a computable model $q_{\Phi}(\tilde{z})$

A Variational Bound

$$\begin{aligned} I(y, \tilde{z}) &= E_{y, \epsilon} \ln \frac{p_{\Phi}(\tilde{z}|y)}{p_{\text{pop}, \Phi}(\tilde{z})} \\ &= E_{y, \epsilon} \ln \frac{p_{\Phi}(\tilde{z}|y)}{q_{\Phi}(\tilde{z})} + E_{y, \epsilon} \ln \frac{q_{\Phi}(\tilde{z})}{p_{\text{pop}, \Phi}(\tilde{z})} \\ &= E_{y, \epsilon} \ln \frac{p_{\Phi}(\tilde{z}|y)}{q_{\Phi}(\tilde{z})} - KL(p_{\text{pop}, \Phi}(\tilde{z}), q_{\Phi}(\tilde{z})) \\ &\leq E_{y, \epsilon} \ln \frac{p_{\Phi}(\tilde{z}|y)}{q_{\Phi}(\tilde{z})} \end{aligned}$$

A Fundamental Equation for the Continuous Case

$$\tilde{z} = z_{\Phi}(y) + f_{\Phi}(y, \epsilon)$$

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y, \epsilon} \ln \frac{p_{\Phi}(\tilde{z}|y)}{q_{\Phi}(\tilde{z})} + \lambda \operatorname{Dist}(y, y_{\Phi}(\tilde{z}))$$

END