

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2020

Noisy Channel RDAs

Rate-Distortion Autoencoders (RDAs)

We compress a continuous signal y to a bit string $\tilde{z}_\Phi(y)$.

We decompress $\tilde{z}_\Phi(y)$ to $y_\Phi(\tilde{z}_\Phi(y))$.

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim P_{\text{op}}} \left[|\tilde{z}_\Phi(y)| + \lambda \text{Dist}(y, y_\Phi(\tilde{z}_\Phi(y))) \right]$$

Rate-Distortion Autoencoders (RDAs)

Since rounding is not differentiable we train by replace rounding by additive noise.

$$\mathcal{L}(\Phi) = E_{y \sim P_{\text{op}}} E_{\epsilon} \left\{ \begin{array}{l} -\ln p_{\Phi}(z_{\Phi}(y) + \epsilon) \\ + \lambda \text{Dist}(y, y_{\Phi}(z_{\Phi}(y) + \epsilon)) \end{array} \right.$$

A noisy-channel RDA uses the noise version without rounding.

Noisy Channel RDAs

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

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

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

Noisy Channel RDAs

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

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

Using parameter-independent noise is called the “reparameterization trick” and allows SGD.

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

Mutual Information as a Channel Rate

Typically $z_{\Phi}(y, \epsilon)$ is simple. For example

$$\epsilon \sim \mathcal{N}(0, I)$$

$$z_{\Phi}(y, \epsilon) = \mu_{\Phi}(y) + \sigma_{\Phi}(y) \odot \epsilon$$

In this example $p_{\Phi}(z|y)$ is easily computed.

Mutual Information Replaces Rate

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

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

A Variational Bound

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

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

Instead we will use a model $\hat{p}_{\Phi}(z)$ to approximate $p_{\text{pop},\Phi}(z)$.

A Variational Bound

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

The Noisy Channel RDA

General Noisy Channel RDA:

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y,\epsilon} \ln \frac{p_{\Phi}(z_{\Phi}(y, \epsilon)|y)}{\hat{p}_{\Phi}(z_{\Phi}(y, \epsilon))} + \lambda \text{Dist}(y, y_{\Phi}(z_{\Phi}(y, \epsilon)))$$

Uniform Box Noise (Rounding) RDA:

$$\begin{aligned} \Phi^* = \operatorname{argmin}_{\Phi} E_y E_{\epsilon \sim [-1/2, 1/2]^d} \\ - \ln \hat{p}_{\Phi}(z_{\Phi}(y) + \epsilon) + \lambda \text{Dist}(y, y_{\Phi}(z_{\Phi}(y) + \epsilon)) \end{aligned}$$

END