

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2020

Noise Contrastive Estimation

Noise Contrastive Estimation

Gutmann and Hyvärinen, 2010

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} - \ln P_{\Psi}(i | y_1, \dots, y_N)$$

p_{Φ} is fixed “noise”

Assume p_{Φ} is both samplable and computable — we can sample from p_{Φ} and for any given y we can compute $p_{\Phi}(y)$.

Assume $P_{\Psi}(i | y_1, \dots, y_N) = \operatorname{softmax}_i s_{\Psi}(y_i)$

Assume Ψ universal

Noise Contrastive Estimation

$$\Psi^* = \underset{\Psi}{\operatorname{argmin}} E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} - \ln P_{\Psi}(i | y_1, \dots, y_N)$$

p_{Φ} is fixed “noise”

Theorem: $\operatorname{pop}(y) = \operatorname{softmax}_y \left(s_{\Psi^*}(y) + \ln p_{\Phi}(y) \right)$

We then have a computable score function (energy function) for the population. We do not have the partition function Z .

Noise Contrastive Estimation

$$\Psi^* = \underset{\Psi}{\operatorname{argmin}} E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} - \ln P_{\Psi}(i | y_1, \dots, y_N)$$

p_{Φ} is fixed “noise”

Lemma: $P_{\Psi^*}(i | y_1, \dots, y_N) = \operatorname{softmax}_i \ln \frac{\operatorname{pop}(y_i)}{p_{\Phi}(y_i)}$

Lemma Proof

$$\begin{aligned}\tilde{p}_{\Phi}^{(N)}(i \text{ and } y_1, \dots, y_N) &= \frac{1}{N} \text{pop}(y_i) \prod_{j \neq i} p_{\Phi}(y_j) \\ &= \alpha \frac{\text{pop}(y_i)}{p_{\Phi}(y_i)}, \quad \alpha = \frac{1}{N} \prod_i p_{\Phi}(y_i)\end{aligned}$$

$$\begin{aligned}\tilde{p}_{\Phi}^{(N)}(i \mid y_1, \dots, y_N) &= \frac{\tilde{p}_{\Phi}^{(N)}(i \text{ and } y_1, \dots, y_N)}{\sum_i \tilde{p}_{\Phi}^{(N)}(i \text{ and } y_1, \dots, y_N)} = \frac{1}{Z} \frac{\text{pop}(y_i)}{p_{\Phi}(y_i)} \\ &= \text{softmax}_i \left(\ln \frac{\text{pop}(y_i)}{p_{\Phi}(y_i)} \right)\end{aligned}$$

Theorem Proof

$$\operatorname{softmax}_i s_{\Psi^*}(y_i) = \operatorname{softmax}_i \ln \frac{\operatorname{pop}(y_i)}{p_{\Phi}(y_i)}$$

is solved by

$$s_{\Psi^*}(y) = \ln \frac{\operatorname{pop}(y)}{p_{\Phi}(y)} + \ln Z$$

giving

$$\operatorname{pop}(y) = \frac{1}{Z} \exp(s_{\Psi}(y) + \ln p_{\Phi}(y))$$

Another Theorem

$$\begin{aligned} E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} & - \ln p_{\Psi^*}^{(N)}(i | y_1, \dots, y_N) \\ & \geq \ln N - \frac{N-1}{N} (KL(\text{pop}, p_{\Phi}) + KL(p_{\Phi}, \text{pop})) \end{aligned}$$

Note that the bound holds with equality for $p_{\Phi} = \text{pop}$.

This is analogous to the JSD expression for the optimal discriminator.

Proof Part A.

$$\begin{aligned} & E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} \ln p_{\Psi^*}(i | y_1, \dots, y_N) \\ &= E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} \ln \left(\operatorname{softmax}_i \ln \frac{\operatorname{pop}(y_i)}{p_{\Phi}(y_i)} \right) [i] \\ &= E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} \ln \frac{\operatorname{pop}(y_i)}{p_{\Phi}(y_i)} - \ln \left(\sum_j \frac{\operatorname{pop}(y_j)}{p_{\Phi}(y_j)} \right) \\ &= E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} \ln \frac{\operatorname{pop}(y_i)}{p_{\Phi}(y_i)} - \ln \left(\frac{1}{N} \sum_j \frac{\operatorname{pop}(y_j)}{p_{\Phi}(y_j)} \right) - \ln N \end{aligned}$$

Proof Part B.

$$\begin{aligned}
& E_{(i, y_1, \dots, y_N) \sim \tilde{p}_\Phi^{(N)}} \ln \frac{\text{pop}(y_i)}{p_\Phi(y_i)} - \ln \left(\frac{1}{N} \sum_j \frac{\text{pop}(y_j)}{p_\Phi(y_j)} \right) - \ln N \\
& \leq E_{(i, y_1, \dots, y_N) \sim \tilde{p}_\Phi^{(N)}} \ln \frac{\text{pop}(y_i)}{p_\Phi(y_i)} - \frac{1}{N} \sum_j \ln \frac{\text{pop}(y_j)}{p_\Phi(y_j)} - \ln N \\
& = E_{y \sim \text{pop}} \ln \frac{\text{pop}(y)}{p_\Phi(y)} - E_{(i, y_1, \dots, y_N) \sim \tilde{p}_\Phi^{(N)}} \frac{1}{N} \sum_j \ln \frac{\text{pop}(y_j)}{p_\Phi(y_j)} - \ln N \\
& = \frac{N-1}{N} (KL(\text{pop}, p_\Phi) + KL(p_\Phi, \text{pop})) - \ln N
\end{aligned}$$

Noise Discriminative Estimation

As in noise contrastive estimation we assume a noise distribution p_Φ that is both samplable and computable.

For noise discriminative estimation, and for $i \in \{-1, 1\}$, we define a probability distribution over pairs $\langle i, y \rangle$ as in GANs.

$$\begin{aligned}\tilde{p}_\Phi(i = 1) &= 1/2 \\ \tilde{p}_\Phi(y|i = 1) &= \text{pop}(y) \\ \tilde{p}_\Phi(y|i = -1) &= p_\Phi(y)\end{aligned}$$

Noise Discriminative Estimation

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{\langle i, y \rangle \sim \tilde{p}_{\Phi}} - \ln P_{\Psi}(i|y)$$

$$\text{Assume } P_{\Psi}(i|y) = \operatorname{softmax}_i \quad is_{\Psi}(y) = \frac{1}{1+e^{-2is_{\Psi}(y)}}$$

Assume Ψ universal.

Noise Discriminative Estimation

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i,y) \sim \tilde{p}_{\Phi}} - \ln P_{\Psi}(i|y)$$

Theorem: $\text{pop}(y) = \text{softmax}_y \quad s_{\Psi^*}(y) + \ln p_{\Phi}(y)$

As with noise contrastive estimation, we have a computable score function (energy function) for the population. We do not have the partition function Z .

Noise Discriminative Estimation

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i, y_1, \dots, y_N) \sim \tilde{p}_{\Phi}^{(N)}} - \ln P_{\Psi}(i | y_1, \dots, y_N)$$

Lemma: $P_{\Psi^*}(i | y) = \operatorname{softmax}_i s_{\Psi^*}(y)$

Lemma Proof

$$\begin{aligned}\tilde{p}_{\Phi}^{(N)}(i \text{ and } y_1, \dots, y_N) &= \frac{1}{N} \text{pop}(y_i) \prod_{j \neq i} p_{\Phi}(y_j) \\ &= \alpha \frac{\text{pop}(y_i)}{p_{\Phi}(y_i)}, \quad \alpha = \frac{1}{N} \prod_i p_{\Phi}(y_i)\end{aligned}$$

$$\begin{aligned}\tilde{p}_{\Phi}^{(N)}(i \mid y_1, \dots, y_N) &= \frac{\tilde{p}_{\Phi}^{(N)}(i \text{ and } y_1, \dots, y_N)}{\sum_i \tilde{p}_{\Phi}^{(N)}(i \text{ and } y_1, \dots, y_N)} = \frac{1}{Z} \frac{\text{pop}(y_i)}{p_{\Phi}(y_i)} \\ &= \text{softmax}_i \left(\ln \frac{\text{pop}(y_i)}{p_{\Phi}(y_i)} \right)\end{aligned}$$

Theorem Proof

$$\operatorname{softmax}_i s_{\Psi^*}(y_i) = \operatorname{softmax}_i \ln \frac{\operatorname{pop}(y_i)}{p_{\Phi}(y_i)}$$

is solved by

$$s_{\Psi^*}(y) = \ln \frac{\operatorname{pop}(y)}{p_{\Phi}(y)} + \ln Z$$

giving

$$\operatorname{pop}(y) = \frac{1}{Z} \exp(s_{\Psi}(y) + \ln p_{\Phi}(y))$$

END