

TTIC 31230 Fundamentals of Deep Learning
Quiz 1

Problem 1: Optimizing Cross Entropy. For this problem we consider a population distribution Pop on the non-negative natural numbers $k \geq 0$. We will work with the population mean $\mu = E_{k \sim \text{Pop}} k$. We consider model distributions defined by the single parameter λ with $0 \leq \lambda < 1$ and defined by the distribution

$$Q_\lambda(k) = (1 - \lambda)\lambda^k$$

(a) Given an expression for the cross-entropy $H(\text{Pop}, Q_\lambda)$ in terms of μ and λ .

Solution: By definition,

$$\begin{aligned} H(\text{Pop}, Q_\lambda) &= -\mathbb{E}_{k \sim \text{Pop}}[\ln Q_\lambda(k)] \\ &= -\mathbb{E}_{k \sim \text{Pop}}[\ln(1 - \lambda) + k \ln \lambda] \\ &= -(\ln(1 - \lambda) + \mu \ln \lambda) . \end{aligned}$$

(b) Solve for the optimal value λ^* minimizing $H(\text{Pop}, Q_\lambda)$ as a function of μ .

Solution:

$$\partial_\lambda H(\text{Pop}, Q_\lambda) = \frac{1}{1 - \lambda} - \frac{\mu}{\lambda} .$$

Setting it to zero and solving for λ yields $\lambda^* = \frac{\mu}{\mu + 1}$.

To see that this is indeed a minimizer, check that

$$\partial_\lambda^2 H(\text{Pop}, Q_\lambda) = \frac{1}{(1 - \lambda)^2} + \frac{\mu}{\lambda^2} > 0$$

for all $\lambda \in [0, 1), \mu \in [0, \infty)$, hence $H(\text{Pop}, Q_\lambda)$ is strictly convex in λ .

(c) Solve for mean value of the distribution Q_{λ^*} in terms of μ .

Solution:

Knowing that Q_λ is a geometric distribution with $p = 1 - \lambda$, it follows that $\mathbb{E}_{k \sim Q_\lambda}[k] = \frac{1-p}{p} = \frac{\lambda}{1-\lambda}$. Plugging $\lambda = \frac{\mu}{\mu+1}$ yields that the mean is μ .

Alternatively, we can avoid using the above fact as follows

$$\begin{aligned}
\mathbb{E}_{k \sim Q_\lambda}[k] &= \sum_{k=0}^{\infty} k(1-\lambda)\lambda^k = (1-\lambda) \sum_{k=0}^{\infty} \lambda k \lambda^{k-1} = (1-\lambda) \sum_{k=0}^{\infty} \lambda \cdot \partial_\lambda(\lambda^k) \\
&= (1-\lambda)\lambda \cdot \partial_\lambda \left(\sum_{k=0}^{\infty} \lambda^k \right) = (1-\lambda)\lambda \cdot \partial_\lambda \left(\frac{1}{1-\lambda} \right) = (1-\lambda)\lambda \frac{1}{(1-\lambda)^2} \\
&= \frac{\lambda}{1-\lambda},
\end{aligned}$$

where we used the fact that the series $\sum_{k=0}^{\infty} \lambda^k$ converges uniformly to $\frac{1}{1-\lambda}$.

Problem 2. Maximum Mutual Information Training. Consider a population distribution Pop on pairs $\langle x, y \rangle$ and a model distribution $Q_\Phi(\hat{y}|x)$. Consider a distribution P_Φ on triples x, y, \hat{y} where $\langle x, y \rangle$ is drawn from Pop and \hat{y} is drawn from $Q_\Phi(\hat{y}|x)$. Under the distribution P_Φ the mutual information between y and \hat{y} is defined by

$$\begin{aligned}
I_\Phi(y, \hat{y}) &= KL(P_\Phi(y, \hat{y}), \text{Pop}(y)P_\Phi(\hat{y})) \\
P_\Phi(y, \hat{y}) &= \sum_x \text{Pop}(x) \text{Pop}(y|x) Q_\Phi(\hat{y}|x) \\
&= E_{x \sim \text{Pop}} \text{Pop}(y|x) Q_\Phi(\hat{y}|x) \\
\text{Pop}(y) &= E_{x \sim \text{Pop}} \text{Pop}(y|x) \\
P_\Phi(\hat{y}) &= E_{x \sim \text{Pop}} Q_\Phi(\hat{y}|x)
\end{aligned}$$

Here we are interested in comparing the fundamental cross entropy objective to the objective of maximizing the mutual information $I_\Phi(y, \hat{y})$.

$$\begin{aligned}
\Phi_1^* &= \underset{\Phi}{\operatorname{argmin}} E_{\langle x, y \rangle \sim \text{Pop}} - \ln Q_\Phi(y|x) \\
\Phi_2^* &= \underset{\Phi}{\operatorname{argmax}} I_\Phi(y, \hat{y})
\end{aligned}$$

(a) Suppose that there exists a perfect predictor – a parameter setting Φ^* such that $Q_{\Phi^*}(\hat{y}|x) = 1$ for $\hat{y} = y$ and zero otherwise. Show using an explicit calculation and standard information theoretic inequalities that a perfect predictor is an optimum of both the cross-entropy objective and the maximum mutual information objective.

Solution:

Since $Q_{\Phi^*}(\hat{y}|x) = 1$ for $\hat{y} = y$ and 0 otherwise, we get

$$-\mathbb{E}_{x, y \sim \text{Pop}}[\ln Q_{\Phi^*}(y|x)] = -\mathbb{E}_{x, y \sim \text{Pop}}[\ln 1] = 0,$$

showing that Φ^* minimizes the cross-entropy as it is lower-bounded by 0.

The above implies $Q_{\Phi^*}(y|x) = \text{Pop}(y|x)$ so $P_{\Phi^*}(y, \hat{y}) = \sum_x \text{Pop}(x) Q_{\Phi^*}(y|x) Q_{\Phi^*}(\hat{y}|x)$, therefore we get that $P_{\Phi^*}(y, y) = \text{Pop}(y)$ for any y , and $P_{\Phi^*}(y, \hat{y}) = 0$ whenever $y \neq \hat{y}$. With this,

$$\begin{aligned} I(y, \hat{y}) &= \sum_y \sum_{y'} P_{\Phi^*}(y, \hat{y}) \log \frac{P_{\Phi^*}(y, \hat{y})}{\text{Pop}(y) P_{\Phi^*}(\hat{y})} \\ &= \sum_y \text{Pop}(y) \log \frac{\text{Pop}(y)}{\text{Pop}(y) \text{Pop}(y)} \\ &= H(y) \end{aligned}$$

Since $I(y, \hat{y}) \leq \min(H(y), H(\hat{y}))$, it follows that Φ^* is a maximizer of the mutual information objective.

(b) Consider binary classification where we have $y, \hat{y} \in \{-1, 1\}$. Show using an explicit calculation and standard information-theoretic inequalities that a perfect anti-predictor with $Q_{\Phi}(\hat{y}|x) = 1$ for $\hat{y} = -y$ is also optimal for the maximum mutual information objective.

Solution: In this case we have that $Q_{\Phi}(y|x) = 1 - \text{Pop}(y|x)$, therefore for any $y \in \{\pm 1\}$ we have $P_{\Phi^*}(y, -y) = \text{Pop}(y)$ and $P_{\Phi^*}(y, y) = 0$. Then, we get

$$\begin{aligned} I(y, \hat{y}) &= \sum_y \sum_{y'} P_{\Phi^*}(y, \hat{y}) \log \frac{P_{\Phi^*}(y, \hat{y})}{\text{Pop}(y) P_{\Phi^*}(\hat{y})} \\ &= \sum_y P_{\Phi^*}(y, -y) \log \frac{P_{\Phi^*}(y, -y)}{\text{Pop}(y) P_{\Phi^*}(-y)} \\ &= \sum_y \text{Pop}(y) \log \frac{\text{Pop}(y)}{\text{Pop}(y) \text{Pop}(y)} \\ &= H(y), \end{aligned}$$

therefore Φ^* is also a maximizer of the MI objective.

Problem 3. Backpropagation for Layer Normalization. Layer normalization is an alternative to batch normalization and is used in the transformer to handle “covariate shift”. In the transformer each a layer has positions in the text that I will index by t and neurons at each position that I will index by i . We can think of this as a sequence of vectors $L[t, I]$. Layer normalization is defined by the following equations where the vectors $A_{\ell+1}[I]$ and $B_{\ell+1}[I]$ are

trained parameters and σ is an arbitrary activation function, typically ReLU.

$$\begin{aligned}\mu_\ell &= \frac{1}{TI} \sum_{t,i} L_\ell[t, i] \\ \sigma_\ell &= \sqrt{\frac{1}{TI} \sum_{t,i} (L_\ell[t, i] - \mu_\ell)^2} \\ \tilde{L}_{\ell+1}[t, i] &= \sigma \left(\frac{A_{\ell+1}[i]}{\sigma_\ell} (L_\ell[t, i] - \mu_\ell) + B_{\ell+1}[i] \right)\end{aligned}$$

Write backpropagation equations for these three assignments.

Solution:

Let $h[t, i] = \frac{L[t, i] - \mu}{\sigma}$

$$\begin{aligned}A.grad[i] &+= \sum_t \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) \frac{L[t, i] - \mu}{\sigma} \\ B.grad[i] &+= \sum_t \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) \\ \sigma.grad &= \sum_{t,i} \tilde{L}.grad[t, i] \sigma'(h[t, i]) \frac{A[i](L[t, i] - \mu)}{\sigma^2} \\ \mu.grad &= \sum_{t,i} \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) \frac{A[i]}{\sigma} \\ L.grad[t, i] &+= \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) \frac{A[i]}{\sigma} \\ \mu.grad &+= \sum_{t,i} \sigma.grad \frac{1}{TI\sigma} (\mu - L[t, i]) = 0 \\ L.grad[t, i] &+= \sigma.grad \frac{1}{TI\sigma} (L[t, i] - \mu) \\ L.grad[t, i] &+= \mu.grad \frac{1}{TI}\end{aligned}$$

Combining the above, we get:

$$\begin{aligned}A.grad[i] &+= \sum_t \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) h[t, i] \\ B.grad[i] &+= \sum_t \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) \\ L.grad[t, i] &+= \frac{1}{\sigma} \left(\tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) A[i] - \frac{\sum_{t,i} \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) A[i]}{TI} \right. \\ &\quad \left. - h[t, i] \frac{\sum_{t,i} \tilde{L}.grad[t, i] \sigma'(A[i]h[t, i] + B[i]) A[i] h[t, i]}{TI} \right)\end{aligned}$$