Assume that probability distributions $P(y)$ are discrete with $\sum_y P(y) = 1$.

**Problem 1:** The problem of population density estimation is defined by the following equation.

$$\Phi^* = \underset{\Phi}{\text{argmin}} \ H(\text{Pop}, P_\Phi) = E_{y \sim \text{Pop}} \ -\log \ P_\Phi(y)$$

This equation is used for language modeling — estimating the probability distribution over the population of English sentences that appear, say, in the New York Times. Show the following.

$$\Phi^* = \underset{\Phi}{\text{argmin}} \ H(\text{Pop}, P_\Phi) = \underset{\Phi}{\text{argmin}} \ KL(\text{Pop}, P_\Phi)$$

Assuming that the model probability $P_\Phi(y)$ can be computed for any given $y$, but that we have no way of computing $\text{Pop}(y)$ for a given $y$, explain why gradient descent on the cross-entropy objective can be done while gradient descent on the KL-divergence form is problematic.

**Problem 2:** Consider the objective

$$P^* = \underset{P}{\text{argmin}} \ H(P, Q) \tag{1}$$

Define $y^*$ by

$$y^* = \underset{y}{\text{argmax}} \ Q(y)$$

Let $\delta_y$ be the distribution such that $\delta_y(y) = 1$ and $\delta_y(y') = 0$ for $y' \neq y$. Show that $\delta_{y^*}$ minimizes (1).
Next consider

$$P^* = \underset{P}{\text{argmin}} \ KL(P, Q) \tag{2}$$

Show that $Q$ is the minimizer of (2).
Next consider a subset $S$ of the possible values and let $Q_S$ be the restriction of $Q$ to the set $S$.

$$Q_S(y) \quad = \quad \frac{1}{Q(S)} \begin{cases} Q(y) & \text{for } y \in S \\ 0 & \text{otherwise} \end{cases}$$

Show that that $KL(Q_S, Q) = -\ln Q(S)$, which will be quite small if $S$ covers much of the mass. Show that, in contrast, $KL(Q, Q_S)$ is infinite unless $Q_S = Q$.
When we optimize a model $P_\Phi$ under the objective $KL(P_\Phi, Q)$ we can get that $P_\Phi$ covers only one high probability region (a mode) of $Q$ (a problem called mode collapse) while optimizing $P_\Phi$ under the objective $KL(Q, P_\Phi)$ we will tend to

get that $P_\Phi$ covers all of $Q$. The two directions are very different even though both are minimized at $P = Q$.

**Problem 5.** Prove the data processing inequality that for any function $f$ with $z = f(y)$ we have $H(z) \leq H(y)$.

**Problem 3:** Consider a joint distribution $P(x, y)$ on discrete random variables $x$ and $y$. We define the marginal distributions $P(x)$ and $P(y)$ as follows.

$$P(x) \quad = \quad \sum_y P(x, y)$$

$$P(y) \quad = \quad \sum_x P(x, y)$$

Let $Q(x, y)$ be defined to be the product of marginals.

$$Q(x, y) = P(x)P(y).$$

We define mutual information by

$$I(x, y) = KL(P, Q)$$

which I will write as

$$I(x, y) = KL(P(x, y), Q(x, y))$$

We define conditional entropy $H(y|x)$ by

$$H(y|x) = E_{x,y} \ -\log P(y|x).$$

(a) Show
$$I(x, y) = H(y) - H(y|x) = H(x) - H(x|y)$$

(b) Explain why (a) implies $H(x) \geq H(x|y)$.

(c) By stating (b) conditioned on $z$ we have

$$H(x|z) \geq H(x|y, z).$$

Use this to show that the data process inequality applies to mutual information, i.e., that for $z = f(y)$ we have $I(x, z) \leq I(x, y)$.

**Problem 4:** (a) For three distributions $P$, $Q$ and $G$ show the following equality.

$$KL(P, Q) = \left( E_{y \sim P} \ \log \frac{G(y)}{Q(y)} \right) + KL(P, G)$$

(b) Show that this implies

$$KL(P,Q) = \sup_G \; E_{y\sim P} \; \log \frac{G(y)}{Q(y)}$$

(c) Now define

$$G(y) \;\; = \;\; \frac{1}{Z} \, Q(y) e^{s(y)}$$

$$Z \;\; = \;\; \sum_y Q(y) e^{s(y)}$$

Show that a distribution $G(y)$ which does not assign zero to any point can be represented by a score $s(y)$ and that under this change of variables we have

$$KL(P,Q) = \sup_s \; E_{y\sim P} \; s(y) - \log E_{y\sim Q} \; e^{s(y)}$$

This is the Donsker-Varadhan variational representation of KL-divergence. This can be used in cases where we can sample from $P$ and $Q$ but cannot compute $P(y)$ or $Q(y)$. Instead we can use a model score $s_\Phi(y)$ where $s_\Phi(y)$ can be computed.