

WRITTEN EXERCISES

1. Maximum Likelihood and KL Divergence. In machine learning, we often need to assess the similarity between pairs of distributions. This is often done using the [Kullback-Leibler](#) (KL) divergence:

$$KL(p(x)||q(x)) = \mathbb{E}_{p(x)} [\log p(x) - \log q(x)]$$

The KL divergence is always non-negative, and equals zero when p and q are identical. This makes it a natural tool for comparing distributions.

This question explores connections between the KL divergence and maximum likelihood learning. Suppose we want to learn a supervised probabilistic model $p_\theta(y|x)$ (e.g., logistic regression) with parameters θ over a dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)}) \mid i = 1, 2, \dots, n\}$. Let $\hat{p}(x, y)$ denote the *empirical* distribution of the data, which is the distribution that assigns a probability of $1/n$ to each of the data points in \mathcal{D} (and zero to all the other possible (x, y)):

$$\hat{p}(x, y) = \begin{cases} \frac{1}{n} & \text{if } (x, y) \in \mathcal{D} \\ 0 & \text{otherwise.} \end{cases}$$

The empirical distribution can be seen as a guess of the true data distribution from which the dataset \mathcal{D} was sampled i.i.d.: it assigns a uniform probability to every possible training instance seen so far, and does not assign any probability to unseen training instances.

Prove that selecting parameters θ by maximizing the likelihood is equivalent to selecting θ that minimize the average KL divergence between the data distribution and the model distribution:

$$\arg \max_{\theta} \mathbb{E}_{\hat{p}(x, y)} [\log p_\theta(y|x)] = \arg \min_{\theta} \mathbb{E}_{\hat{p}(x)} [KL(\hat{p}(y|x) || p_\theta(y|x))].$$

Here, $\mathbb{E}_{p(x)} f(x)$ denotes $\sum_{x \in \mathcal{X}} f(x)p(x)$ if x is discrete and $\int_{x \in \mathcal{X}} f(x)p(x)dx$ if x is continuous.

$$\arg \max_{\theta} \mathbb{E}_{\hat{p}(x, y)} [\log p_{\theta}(y|x)] = \arg \min_{\theta} \mathbb{E}_{\hat{p}(x)} [KL(\hat{p}(y|x) || p_{\theta}(y|x))]$$

$$\text{RHS} = \arg \min_{\theta} \mathbb{E}_{\hat{p}(x)} [\mathbb{E}_{\hat{p}(y|x)} [\log p(y|x) - \log p_{\theta}(y|x)]]$$

$$= \arg \min_{\theta} \mathbb{E}_{\hat{p}(x)} [\mathbb{E}_{p(y|x)} [\log \hat{p}(y|x) - \log p_{\theta}(y|x)]]$$

$$= \arg \min_{\theta} \mathbb{E}_{\hat{p}(x)} \left[\sum_{(x, y) \in \mathcal{X}} \log \hat{p}(y|x) \cdot \frac{P(x, y)}{P(x)} - \sum_{(x, y) \in \mathcal{X}} \log p_{\theta}(y|x) \cdot \frac{P(x, y)}{P(x)} \right]$$

$$= \arg \min_{\theta} \mathbb{E}_{\hat{p}(x)} \left[\left(\sum_y \log \hat{p}(y|x) - \sum_y \log p_{\theta}(y|x) \right) \cdot \frac{P(x, y)}{P(x)} \right]$$

$$= \arg \min_{\theta} \sum_x \sum_y (\log \hat{p}(y|x) - \log p_{\theta}(y|x)) \cdot \frac{P(x, y)}{P(x)} \cdot P(x)$$

$$= \arg \min_{\theta} \sum_x \sum_y (\log \hat{p}(y|x) - \log p_{\theta}(y|x)) \cdot P(x, y)$$

$$\begin{aligned}
 \because Q \text{ is not dependent on } \log p(y|x) \\
 \therefore \text{LHS} &= \arg\min_a \sum_x \sum_y (-\log p_a(y|x) \cdot p(x,y)) \\
 &= \arg\min_a E_{p(x,y)} (-\log p_a(y|x)) \\
 &= \arg\max_a E_{p(x,y)} [\log p_a(y|x)]
 \end{aligned}$$

$$\therefore \text{RHS} = \text{LHS}$$

2. Gradient and log-likelihood for logistic regression.

(a) Let $\sigma(a) = \frac{1}{1+e^{-a}}$ be the sigmoid function. Show that $\frac{d\sigma(a)}{da} = \sigma(a)(1-\sigma(a))$.

(b) Using the previous result and the chain rule of calculus, derive the expression for the gradient of the log likelihood:

$$\nabla \ell(\theta) = [y - \sigma(\theta^T \mathbf{x})] \mathbf{x}$$

where

$$\ell(\theta) = y \log \sigma(\theta^T \mathbf{x}) + (1-y) \log (1 - \sigma(\theta^T \mathbf{x}))$$

$$\begin{aligned}
 \text{(a). } \frac{d\sigma(a)}{da} &= \frac{d}{da} \left(\frac{1}{1+e^{-a}} \right) \\
 &= \frac{-1}{(1+e^{-a})^2} \cdot e^{-a} = \frac{e^{-a}}{(1+e^{-a})^2}
 \end{aligned}$$

$$\because 1 - \sigma(a) = 1 - \frac{1}{1+e^{-a}} = \frac{1+e^{-a}-1}{1+e^{-a}} = \frac{e^{-a}}{1+e^{-a}}$$

$$\therefore \sigma(a) \cdot (1 - \sigma(a)) = \frac{e^{-a}}{1+e^{-a}} \cdot \frac{1}{1+e^{-a}} = \frac{e^{-a}}{(1+e^{-a})^2} = \frac{d\sigma(a)}{da}$$

$$\text{(b). } \ell(a) = y \log \sigma(a^T \mathbf{x}) + (1-y) \log (1 - \sigma(a^T \mathbf{x}))$$

$$\begin{aligned}
&= y \nabla \log \sigma(\omega^T x) + (1-y) \nabla \log (1-\sigma(\omega^T x)) \\
&= y \cdot \frac{1}{\sigma(\omega^T x)} \cdot \nabla \sigma(\omega^T x) + (1-y) \frac{1}{1-\sigma(\omega^T x)} \cdot \nabla (1-\sigma(\omega^T x)) \\
&= y \cdot \frac{1}{\sigma(\omega^T x)} \cdot x \sigma(\omega^T x) (1-\sigma(\omega^T x)) + (1-y) \cdot \frac{1}{1-\sigma(\omega^T x)} \cdot (-\sigma(\omega^T x) (1-\sigma(\omega^T x))) x \\
&= y \cdot \frac{\sigma(\omega^T x) (1-\sigma(\omega^T x)) x}{\sigma(\omega^T x)} + (1-y) \cdot \frac{-\sigma(\omega^T x) (1-\sigma(\omega^T x)) x}{1-\sigma(\omega^T x)} \\
&= xy - y \cancel{\sigma(\omega^T x)} x - x \cancel{\sigma(\omega^T x)} + y \cancel{\sigma(\omega^T x)} x \\
&= xy - x \sigma(\omega^T x) = x(y - \sigma(\omega^T x))
\end{aligned}$$

3. Analytical solution of the Ordinary Least Squares Estimation. Consider we have a simple dataset of n labeled data $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where data $x_i \in \mathbb{R}$ and $y_i \in \mathbb{R}$ is its corresponding label. We use a simple estimated regression function of:

$$\hat{y}_i = \theta_0 + \theta_1 x_i$$

Instead of gradient descent which works in an iterative manner, we try to directly solve this problem. We define the cost function as the residual sum of squares, parameterized by θ_0, θ_1 :

$$J(\theta_0, \theta_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- (a) Calculate the partial derivative of $\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ and $\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$.
- (b) Consider the fact that $J(\theta_0, \theta_1)$ has a unique optimum, we can actually get the analytical solution of θ_0, θ_1 by the following normal equations:

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = 0$$

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = 0$$

prove the following proprieties that

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

and

$$\theta_1 = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})}$$

(Note: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.)

- (c) Calculate the sum of the residuals $\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))$. What can you learn from the value of $\sum_{i=1}^n e_i$?

$$(a). \frac{\partial}{\partial \alpha_0} J(\alpha_0, \alpha_1) = \frac{\partial}{\partial \alpha_0} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{\partial}{\partial \alpha_0} (y_i - \hat{y}_i)^2$$

$$= 2(y_i - \hat{y}_i) \cdot (y_i - \hat{y}_i)'$$

$$= 2(y_i - \hat{y}_i) \cdot (-\nabla \hat{y}_i)$$

$$= (2\hat{y}_i - 2y_i) \cdot (\nabla(\alpha_0 + \alpha_1 x_i))$$

$$= (2\alpha_0 + 2\alpha_1 x_i - 2y_i) \cdot 1$$

$$= 2\alpha_0 + 2\alpha_1 x_i - 2y_i = -2 \sum_{i=1}^n (y_i - \alpha_0 - \alpha_1 x_i)$$

$$\frac{\partial}{\partial \alpha_1} J(\alpha_0, \alpha_1) = 2(y_i - \hat{y}_i) \cdot (y_i - \hat{y}_i)'$$

$$= (2\hat{y}_i - 2y_i) \cdot (-\nabla \hat{y}_i)$$

$$= (2\hat{y}_i - 2y_i) \cdot (\nabla(\alpha_0 + \alpha_1 x_i))$$

$$= (2\alpha_0 + 2\alpha_1 x_i - 2y_i) \cdot x_i$$

$$= -2x_i \sum_{i=1}^n (y_i - \alpha_0 - \alpha_1 x_i)$$

$$b). \textcircled{1} -2 \sum_{i=1}^n (y_i - \alpha_0 - \alpha_1 x_i) = 0$$

$$\sum_{i=1}^n (y_i - \alpha_0 - \alpha_1 x_i) = 0$$

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \alpha_0 - \sum_{i=1}^n \alpha_1 x_i = 0$$

$$n\bar{y} - n\alpha_0 - n\bar{x}\alpha_1 = 0$$

$$n\alpha_0 = n\bar{y} - n\bar{x}\alpha_1$$

$$\alpha_0 = \bar{y} - \alpha_1 \bar{x}$$

$$\therefore \sum_{i=1}^n \alpha_0 = \sum_{i=1}^n (\bar{y} - \alpha_1 \bar{x})$$

$$\therefore \alpha_0 = \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n \alpha_1 x_i$$

$$\therefore \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad \therefore Q_0 = \bar{y} - Q_1 \bar{x}$$

$$(2) \quad -2x_i \sum_{i=1}^n (y_i - Q_1 x_i - Q_0) = 0$$

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i Q_1 x_i - \sum_{i=1}^n x_i Q_0 = 0$$

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i Q_1 x_i - \sum_{i=1}^n x_i \sum_{j=1}^n (\bar{y}_j - Q_1 \bar{x}) = 0$$

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i (\bar{y} - Q_1 \bar{x}) - \sum_{i=1}^n x_i Q_1 x_i = 0$$

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \bar{y} + \sum_{i=1}^n x_i Q_1 \bar{x} - \sum_{i=1}^n x_i Q_1 x_i = 0$$

$$\Rightarrow \sum_{i=1}^n x_i (y_i - \bar{y}) + \sum_{i=1}^n Q_1 x_i (\bar{x} - x_i) = 0$$

$$\therefore Q_1 = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})}$$

$$(c) \quad \sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - (\bar{y} - Q_1 \bar{x} + Q_1 x_i))$$

$$= \sum_{i=1}^n ((y_i - \bar{y}) + Q_1 (x_i - \bar{x}))$$

$$= \sum_{i=1}^n (y_i - \bar{y}) - Q_1 \sum_{i=1}^n (x_i - \bar{x})$$

$$\therefore Q_1 = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} \Rightarrow \sum_{i=1}^n (y_i - \bar{y}) - \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} \cdot \sum_{i=1}^n (x_i - \bar{x})$$

$$\Rightarrow \sum_{i=1}^n (y_i - \bar{y}) - \sum_{i=1}^n (y_i - \bar{y}) = 0$$

$$\therefore \sum_{i=1}^n e_i = 0$$

$\therefore \sum_{i=1}^n e_i = 0 \Rightarrow$ the sum of residual is 0, means that
given α_0, α_1 the regression function $\hat{y}_i = \alpha_0 + \alpha_1 x_i$ can
perfectly predict the data set, it's a perfect estimation.