

Problem Set #3: Deep Learning & Unsupervised Learning

Problem 1 A simple neural network

Let $X = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ be dataset of m examples with 2 features. That is, $x^{(i)} \in \mathbb{R}^2$. Samples are classified into 2 categorie with labels $y \in \{0, 1\}$, as shown in Figure 1. Want to perform binary classification using a simple neural networks with the architecture shown in Figure 2.

Two features x_1 and x_2 , the three neurons in the hidden layer h_1, h_2, h_3 , and the output neuron as o . Weight from x_i to h_j be $w_{i,j}^{[1]}$ for $i = 1, 2$ and $j = 1, 2, 3$, and weight from h_j to o be $w_j^{[2]}$. Finally, denote intercept weight for h_j as $w_{0,j}^{[1]}$ and the intercept weight for o as $w_0^{[2]}$. Use average squared loss instead of the usual negative log-likelihood:

$$l = \frac{1}{m} \sum_{i=1}^m (o^{(i)} - y^{(i)})^2.$$

(a) Suppose we use sigmoid function as activation function for h_1, h_2, h_3 , and o . We have

$$h_1 = g(w_1^{[1]}x), \quad h_2 = g(w_2^{[1]}x), \quad h_3 = g(w_3^{[1]}x), \quad o = g(w^{[2]}h).$$

Hence,

$$\frac{\partial l}{\partial w_{1,2}^{[1]}} = \frac{1}{m} \sum_{i=1}^m 2(o^{(i)} - y^{(i)})o^{(i)}(1 - o^{(i)})w_2^{[2]}h_2^{(i)}(1 - h_2^{(i)})x_1^{(i)},$$

where $h_2^{(i)} = g(w_{0,2}^{[1]} + w_{1,2}^{[1]}x_1^{(i)} + w_{2,2}^{[1]}x_2^{(i)})$ and g is the sigmoid function. Therefore, the gradient descent update to $w_{1,2}^{[1]}$, assuming learning rate α is

$$w_{1,2}^{[1]} := w_{1,2}^{[1]} - \frac{2\alpha}{m} \sum_{i=1}^m (o^{(i)} - y^{(i)})o^{(i)}(1 - o^{(i)})w_2^{[2]}h_2^{(i)}(1 - h_2^{(i)})x_1^{(i)}$$

where $h_2^{(i)} = g(w_{0,2}^{[1]} + w_{1,2}^{[1]}x_1^{(i)} + w_{2,2}^{[1]}x_2^{(i)})$.

(b) Now, suppose the activation function for h_1, h_2, h_3 , and o is the step function $f(x)$, defined as

$$f(x) = \begin{cases} 1, & (x \geq 0), \\ 0, & (x < 0). \end{cases}$$

Is it possible to have a set of weights that allow the neural network to classify this dataset with 100% accuracy? If so, provide a set of weights by completing `optimal_step_weights` within `src/p01_nn.py` and explain your reasoning for those weights. If not, please explain the reasoning.

There is a set of weights that allow the neural network to classify this dataset with 100% accuracy. For the step function activation, we have

$$\begin{aligned} h_1 &= f(w_1^{[1]}x) = f(w_{0,1}^{[1]} + w_{1,1}^{[1]}x_1 + w_{2,1}^{[1]}x_2) \\ h_2 &= f(w_2^{[1]}x) = f(w_{0,2}^{[1]} + w_{1,2}^{[1]}x_1 + w_{2,2}^{[1]}x_2) \\ h_3 &= f(w_3^{[1]}x) = f(w_{0,3}^{[1]} + w_{1,3}^{[1]}x_1 + w_{2,3}^{[1]}x_2) \\ o &= f(w^{[2]}h) = f(w_0^{[2]} + w_1^{[2]}h_1 + w_2^{[2]}h_2 + w_3^{[2]}h_3). \end{aligned}$$

Notice from Figure 1 that the label $y^{(i)} = 0$ if and only if $x^{(i)}$ satisfies

$$\begin{cases} x_2^{(i)} > 0.5, \\ x_1^{(i)} > 0.5, \\ x_1^{(i)} + x_2^{(i)} < 4. \end{cases}$$

Now, let

$$w_1^{[1]} = \begin{bmatrix} 0.5 \\ 0 \\ -1 \end{bmatrix}, \quad w_2^{[1]} = \begin{bmatrix} 0.5 \\ -1 \\ 0 \end{bmatrix}, \quad w_3^{[1]} = \begin{bmatrix} -4 \\ 1 \\ 1 \end{bmatrix}, \quad w^{[2]} = \begin{bmatrix} -0.5 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

Under this set of weights, if all inequalities are satisfied, then $h_1 = h_2 = h_3 = 0$ and $w^{[2]}h = -0.5$. Otherwise, $h_1 + h_2 + h_3 \geq 1$ and $w^{[2]}h \geq -0.5$. Hence, This set of weights will capture all the conditions and allow the neural network to classify this dataset with 100% accuracy.

- (c) Let the activation function for h_1, h_2, h_3 , and o is the linear function $f(x) = x$, and the activation function for o be the same step function as before. Is it possible to have a set of weights that allow the neural network to classify this dataset with 100% accuracy? If so, provide a set of weights by completing `optimal_linear_weights` within `src/p01_nn.py` and explain your reasoning for those weights. If not, please explain the reasoning. ■

Problem 2 KL divergence and maximum likelihood

Kullback-Leibler (KL) divergence is a measure of how much one probability distribution is different from a second one. The *KL divergence* between two discrete-valued distribution $P(X)$, $Q(X)$ over the outcome space \mathcal{X} is defined as follows:

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}.$$

Assume $P(x) > 0$ for all x . (One other standard thing to do is adopt the convention that $0 \log 0 = 0$.) Sometimes, we also write the KL divergence more explicitly as $D_{\text{KL}}(P \parallel Q) = D_{\text{KL}}(P(X) \parallel Q(X))$.

Background on Information Theory

The *entropy* of a probability distribution $P(X)$, defined as

$$H(P) = - \sum_{x \in \mathcal{X}} P(x) \log P(x).$$

measures how dispersed a probability distribution is. Notably, $\mathcal{N}(\mu, \sigma^2)$ has the highest entropy among all possible continuous distribution that has mean μ and variance σ^2 . The entropy $H(P)$ is the best possible long term average bits per message (optimal) that can be achieved under probability distribution $P(X)$.

The *cross entropy* is defined as

$$H(P, Q) = - \sum_{x \in \mathcal{X}} P(x) \log Q(x).$$

The cross entropy $H(P, Q)$ is the long term average bits per message (suboptimal) that results under a distribution $P(X)$, by reusing an encoding scheme designed to be optimal for a scenario with probability distribution $Q(X)$.

Notice that

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log P(x) - \sum_{x \in \mathcal{X}} P(x) \log Q(x) = H(P, Q) - H(P).$$

If $H(P, Q) = 0$, then it necessarily means $P = Q$. In ML, it is common task to find distribution Q that is close to another distribution P . To achieve this, we optimize $D_{\text{KL}}(P \parallel Q)$. Later we will see that Maximum Likelihood Estimation turns out to be equivalent minimizing KL divergence between the training data and the model.

(a) **Nonnegativity.** Prove that

$$D_{\text{KL}}(P \parallel Q) \geq 0$$

and $D_{\text{KL}}(P \parallel Q) = 0$ if and only if $P = Q$.

Hint: Use Jensen's inequality.

Proof. By definition,

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)} = - \sum_{x \in \mathcal{X}} P(x) \log \frac{Q(x)}{P(x)}.$$

Since $-\log x$ is strictly convex, by Jensen's inequality, we have

$$D_{\text{KL}}(P \parallel Q) = - \sum_{x \in \mathcal{X}} P(x) \log \frac{Q(x)}{P(x)} \geq - \log \sum_{x \in \mathcal{X}} P(x) \frac{Q(x)}{P(x)} = 0.$$

When the equality holds,

$$\log \frac{Q(x)}{P(x)} = 0$$

with probability 1. That is, $Q = P$ with probability 1. This completes the proof. \square

- (b) **Chain rule for KL divergence.** The KL divergence between 2 conditional distributions $P(X \mid Y)$, $Q(X \mid Y)$ is defined as follows:

$$D_{\text{KL}}(P(X \mid Y) \parallel Q(X \mid Y)) = \sum_y P(y) \left(\sum_x P(x \mid y) \log \frac{P(x \mid y)}{Q(x \mid y)} \right).$$

This can be thought of as the expected KL divergence between the corresponding conditional distributions on x . That is, between $P(X \mid Y = y)$ and $Q(X \mid Y = y)$, where the expectation is taken over the random y .

Prove the following chain rule for KL divergence:

$$D_{\text{KL}}(P(X, Y) \parallel Q(X, Y)) = D_{\text{KL}}(P(X) \parallel Q(X)) + D_{\text{KL}}(P(Y \mid X) \parallel Q(Y \mid X)).$$

Proof.

$$\begin{aligned} \text{LHS} &= \sum_x \sum_y P(x, y) \log \frac{P(x, y)}{Q(x, y)} \\ &= \sum_x \sum_y P(y \mid x) P(x) \left[\log \frac{P(y \mid x)}{Q(y \mid x)} + \log \frac{P(x)}{Q(x)} \right] \\ &= \sum_x \sum_y P(y \mid x) P(x) \log \frac{P(y \mid x)}{Q(y \mid x)} + \sum_x P(x) \log \frac{P(x)}{Q(x)} \sum_y P(y \mid x) \\ &= \sum_x \sum_y P(y \mid x) P(x) \log \frac{P(y \mid x)}{Q(y \mid x)} + \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &= D_{\text{KL}}(P(X) \parallel Q(X)) + D_{\text{KL}}(P(Y \mid X) \parallel Q(Y \mid X)) \\ &= \text{RHS}. \end{aligned}$$

\square

- (c) **KL and maximum likelihood.** Consider density estimation problem and suppose we are given training set $\{x^{(i)}\}_{i=1}^m$. Let the empirical distribution be $\hat{P}(x) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{x^{(i)} = x\}$. (\hat{P} is just the uniform distribution over the training set; i.e., sampling from the empirical distribution is the same as picking a random example from the training set.)

Suppose we have a family of distributions P_θ parametrized by θ . Prove that finding the maximum likelihood estimates for the parameter θ is equivalent to finding P_θ with minimal KL divergence from \hat{P} . That is, prove that

$$\operatorname{argmin}_{\theta} D_{\text{KL}}(\hat{P} \parallel P_{\theta}) = \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P_{\theta}(x^{(i)}).$$

Proof. Notice that \hat{P} is the uniform distribution over the training set, thus $\hat{P}(x^{(i)}) = \frac{1}{m}$ for $i = 1, \dots, m$. It follows that

$$D_{\text{KL}}(\hat{P} \parallel P_{\theta}) = \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P_{\theta}(x)} = -\log m - \frac{1}{m} \sum_{i=1}^m \log P_{\theta}(x^{(i)}).$$

Hence,

$$\operatorname{argmin}_{\theta} D_{\text{KL}}(\hat{P} \parallel P_{\theta}) = \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P_{\theta}(x^{(i)}),$$

as desired. □

■

Remark: Consider the relationship between parts (b-c) and multi-variate Bernoulli Naive Bayes parameter estimation. In Naive Bayes model we assumed P_{θ} is the following form: $P_{\theta}(x, y) = p(y) \prod_{i=1}^n p(x_i | y)$. By the chain rule for KL divergence, we therefore have

$$D_{\text{KL}}(\hat{P} \parallel P_{\theta}) = D_{\text{KL}}(\hat{P}(y) \parallel p(y)) + \sum_{i=1}^n D_{\text{KL}}(\hat{P}(x_i | y) \parallel p(x_i | y)).$$

This shows that finding the maximum likelihood/minimum KL divergence estimates of the parameters decomposes into $2n + 1$ independent optimization problems: One for the class priors $p(y)$, and one for each conditional distributions $p(x_i | y)$ for each feature x_i given each of the two possible labels for y . Specifically, finding the maximum likelihood estimates for each of these problems individually results in also maximizing the likelihood of the joint distribution. This similarly applies to Bayesian networks. △

Problem 3 KL divergence, Fisher Information, and the Natural Gradient

KL divergence between the two distributions is an asymmetric measure of how different two distributions are. Consider two distributions over the same space given by densities $p(x)$, $q(x)$. The KL divergence between two continuous distributions is defined as

$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= \int p(x) \log \frac{p(x)}{q(x)} dx \\ &= \mathbb{E}_{x \sim p(x)}[\log p(x)] - \mathbb{E}_{x \sim p(x)}[\log q(x)]. \end{aligned}$$

A nice property of KL divergence is that it is invariant to parametrization. This means, KL divergence evaluates to the same value no matter how we parametrize the distribution P and Q . For example, if P and Q are in exponential family, the KL divergence between them is the same whether we are using natural parameters, natural parameters, or canonical parameters, or any arbitrary parametrization.

Now consider the problem of fitting model parameters using gradient descent. While KL divergence is invariant to parametrization, the gradient w.r.t the model parameters gradient is *invariant to parametrization*. We need to use *natural gradient*. This will make the optimization process invariant to the parametrization.

We will construct and derive the natural gradient update rule. Along the way, we will introduce *score function* and *Fisher Information*. Finally, we will see how this new natural gradient based optimization is actually equivalent to Newton's method for Generalized Linear Models.

Let the distribution of a random variable Y parametrized by $\theta \in \mathbb{R}^n$ be $p(y; \theta)$.

- (a) **Score function.** The *score function* with $p(y; \theta)$ is defined as $\nabla_{\theta} \log p(y; \theta)$, which signifies the sensitivity of the likelihood function with respect to the parameters.

Show that the expected value of the score is 0.

Proof. The expected value of the score

$$\begin{aligned} \mathbb{E}_{y \sim p(y; \theta)}[\nabla_{\theta'} \log p(y; \theta)]_{\theta'=\theta} &= \int p(y; \theta) [\nabla_{\theta'} \log p(y; \theta)]_{\theta'=\theta} dy \\ &= \int p(y; \theta) \frac{1}{p(y; \theta)} [\nabla_{\theta'} p(y; \theta)]_{\theta'=\theta} dy \\ &= \left[\nabla_{\theta'} \int p(y; \theta) dy \right]_{\theta'=\theta} \\ &= 0 \end{aligned}$$

□

- (b) **Fisher information.** *Fisher information* is defined as the covariance matrix of the score function,

$$\mathcal{I}(\theta) = \text{cov}_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta')]_{\theta'=\theta}.$$

Intuitively, the Fisher information represents the amount of information that a random variable Y carries about a parameter θ of interest. Show that the Fisher information can be equivalently given by

$$\mathcal{I}(\theta) = \mathbb{E}_{y \sim p(y; \theta)} \left[\left[\nabla_{\theta'} \log p(y; \theta') \nabla_{\theta'} \log p(y; \theta')^T \right]_{\theta' = \theta} \right].$$

Note that the fisher information is a function of the parameter. The parameter is both a) the parameter value at which the score function is evaluated, and b) the parameter of the distribution with respect to which the expectation and variance is calculated.

Proof. Since $\mathbb{E}_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta')]_{\theta' = \theta} = 0$, we have

$$\text{cov}_{y \sim p(y; \theta)} [\nabla_{\theta'} \log p(y; \theta')]_{\theta' = \theta} = \mathbb{E}_{y \sim p(y; \theta)} \left[\left[\nabla_{\theta'} \log p(y; \theta') \nabla_{\theta'} \log p(y; \theta')^T \right]_{\theta' = \theta} \right]$$

by the definition of covariance. This completes the proof. \square

- (c) **Fisher information (alternate form).** It turns out that Fisher information can not only be defined as the covariance of the score function, but in most situations it can also be represented as the expected negative Hessian of the log-likelihood. Show that

$$\mathbb{E}_{y \sim p(y; \theta)} \left[\left[-\nabla_{\theta'}^2 \log p(y; \theta') \right]_{\theta' = \theta} \right] = \mathcal{I}(\theta).$$

Proof. From (b), we know that

$$\begin{aligned} \mathcal{I}_{ij}(\theta) &= \mathbb{E} \left[\frac{\partial}{\partial \theta_i} \log p(y; \theta') \frac{\partial}{\partial \theta_j} \log p(y; \theta') \right] \\ &= \mathbb{E} \left[\frac{1}{p(y; \theta')} \frac{\partial}{\partial \theta_i} p(y; \theta') \frac{\partial}{\partial \theta_j} p(y; \theta') \right]. \end{aligned}$$

Also, for the left hand side of the expression we need to prove, we have

$$\begin{aligned} \text{LHS} &= \mathbb{E} \left[-\frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \nabla_{\theta'}^2 \log p(y; \theta') \right] \\ &= \mathbb{E} \left[\frac{\partial}{\partial \theta_i} \frac{1}{p(y; \theta')} \frac{\partial}{\partial \theta_j} p(y; \theta') \right] \\ &= \mathbb{E} \left[-\frac{1}{p(y; \theta')} \frac{\partial}{\partial \theta_i} p(y; \theta') \frac{\partial}{\partial \theta_j} p(y; \theta') + \frac{1}{p(y; \theta')} \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} p(y; \theta') \right] \end{aligned}$$

For the second term, we have

$$\begin{aligned} \mathbb{E} \left[\frac{1}{p(y; \theta')} \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} p(y; \theta') \right] &= \int \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} p(y; \theta') dy \\ &= \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \int p(y; \theta') dy \\ &= 0. \end{aligned}$$

Hence,

$$\mathbb{E}_{y \sim p(y; \theta)} \left[\left[-\nabla_{\theta'}^2 \log p(y; \theta') \right]_{\theta' = \theta} \right] = \mathbb{E} \left[-\frac{1}{p(y; \theta')} \frac{\partial}{\partial \theta_i} p(y; \theta') \frac{\partial}{\partial \theta_j} p(y; \theta') \right] = \mathcal{I}(\theta).$$

This completes the proof. \square

Remark. This shows that the expected curvature of the log-likelihood function is also equal to the Fisher information matrix. If the curvature of the log-likelihood is steep, this generally means you need fewer number of data samples to estimate that parameter well, and vice versa. The fisher information matrix associated with a statistical model parameterized by θ is extremely important in determining how a model behaves as a function of the number of training set examples. \triangle

- (d) **Approximate D_{KL} with Fisher information.** We are interested in the set of all distributions that are at a small fixed D_{KL} distance away from the current distribution. To calculate KL divergence between $p(y; \theta)$ and $p(y; \theta + d)$, we approximate with Fisher information at θ . Show that

$$D_{\text{KL}}(p_\theta \parallel p_{\theta+d}) \approx \frac{1}{2} d^T \mathcal{I}(\theta) d.$$

Proof. Towards Taylor expansion, we have

$$\begin{aligned} D_{\text{KL}}(p_\theta \parallel p_\theta) &= 0 \\ \nabla_{\theta'} D_{\text{KL}}(p_\theta \parallel p_{\theta'}) &= \nabla_{\theta'} \mathbb{E}[\log p_\theta] - \mathbb{E}[\log p_{\theta'}] = 0 \\ \nabla_{\theta'}^2 D_{\text{KL}}(p_\theta \parallel p_{\theta'}) &= \nabla_{\theta'}^2 \mathbb{E}[\log p_\theta] - \mathbb{E}[\log p_{\theta'}] = \mathcal{I}(\theta). \end{aligned}$$

Hence,

$$D_{\text{KL}}(p_\theta \parallel p_{\hat{\theta}}) \approx \frac{1}{2} d^T \mathcal{I}(\theta) d.$$

\square

- (e) **Natural gradient.** Want to maximize the log-likelihood by moving only by a fixed D_{KL} distance from the current position. Now set up the constrained optimization problem that will yield the natural gradient update d . Let the log-likelihood objective be $\ell(\theta) = \log p(y; \theta)$. Let the D_{KL} distance we want to move by be some small positive constant c . The natural gradient update d^* is

$$d^* = \underset{d}{\operatorname{argmax}} \ell(\theta + d) \text{ subject to } D_{\text{KL}}(p_\theta \parallel p_{\theta+d}) = c.$$

In order to solve this, use Taylor expansion and Lagrangian multipliers.

For the optimization problem, consider the Lagrangian

$$\begin{aligned}
L(d, \lambda) &= \ell(\theta + d) - \lambda(D_{\text{KL}}(p_\theta \parallel p_{\theta+d}) - c) \\
&= \log p(y; \theta + d) - \lambda(D_{\text{KL}}(p_\theta \parallel p_{\theta+d}) - c) \\
&= \log p(y; \theta) + d^T \frac{\nabla_\theta p(y; \theta)}{p(y; \theta)} - \lambda \left(\frac{1}{2} d^T \mathcal{I}(\theta) d - c \right).
\end{aligned}$$

Set

$$\begin{cases} \nabla_d L(d, \lambda) = 0, \\ \nabla_\lambda L(d, \lambda) = 0. \end{cases} \implies \begin{cases} \frac{\nabla_\theta p(y; \theta)}{p(y; \theta)} = \lambda \mathcal{I}(\theta) d, \\ \frac{1}{2} d^T \mathcal{I}(\theta) d = c. \end{cases}$$

This gives

$$d = \sqrt{\frac{2c}{\nabla_\theta p(y; \theta)^T (I^{-1})^T \nabla_\theta p(y; \theta)}} p(y; \theta) I^{-1} \nabla_\theta p(y; \theta).$$

- (f) **Relation to Newton's Method.** Show that the direction of update of Newton's method, and the direction of natural gradient, are exactly the same for GLMs. Refer to results in problem set 1 question 4. For natural gradient, it is sufficient to use \tilde{d} , the unscaled natural gradient

■