

Hw 6

10)

Now, let $x = \log [\pi(y)/q(y)]$ Then we have:

$$E[x] = \int \log [\pi(y)/q(y)] \cdot P(y) dy$$

where $P(y)$ is probability density fn.

Using Jensen's inequality with convex fn.

$$\phi(a) = e^a$$

$$\mathbb{E}[x] \leq E[e^x]$$

Substituting $x = \log [\pi(y)/q(y)]$

$$\mathbb{E}[\log (\pi(y)/q(y))] \leq E[e^{\log [\pi(y)/q(y)]}]$$

Simplifying right hand side,

$$E[e^{\log [\pi(y)/q(y)]}] = E[\pi(y)/q(y)] = \int \pi(y)/q(y) \cdot P(y) dy$$

Since $\log(\cdot)$ and e^{\cdot} are inverse functions, we have

$e^{\log(a)} = a$, substituting this in left hand side & simplifying

$$\pi(y)/q(y) \leq E[\pi(y)/q(y)]$$

Taking natural logarithm of both sides.

$$\log [\pi(y)/q(y)] \leq \log [E[\pi(y)/q(y)]]$$

∴ we shown that

$$\log [\pi(y)/q(y)] \leq \log E[\log [\pi(y)/q(y)]]$$

(Q2)

$$\hat{P}_{\text{Pr}}(U_k(u)) = \frac{1}{(M-m+1)} \sum_{t=m}^M \text{Pr}(u | U_l^{(t)}, l \neq k)$$

To justify the estimate given by eq (8.50), first we need to understand the context of problem. we are given a statistical problem where we have training data $U_l^{(t)}$, $l \neq k$ with l representing class labels and t representing data points. we are interested in estimating Probability $\text{Pr}(u | U_l^{(t)}, l \neq k)$ for a given class u and some data points $U_l^{(t)}$.

we can use relationship provided in eq(1) to compute the probability $\text{Pr}(u | U_l^{(t)}, l \neq k)$ by marginalizing over distribution of data points $U_l^{(t)}$

$$\text{Pr}(u | U_l^{(t)}, l \neq k) = \int \text{Pr}(u | U_l^{(t)}, B, l \neq k) d(\text{Pr}(B))$$

Now, let's divide range of data points into $M-m+1$ equally spaced intervals and let the midpoints of these intervals be denoted as m_1^m, m_2^m, \dots we can approximate the integral by summing midpoints.

$$\text{Pr}(u | U_l^{(t)}, l \neq k) \approx \frac{1}{M-m+1} \sum_{t=m}^M \text{Pr}(u | U_l^{(t)}, B, l \neq k) d(\text{Pr}(B))$$

now as summation does not depend on the measure $d(\text{Pr}(B))$, we can remove it from sum:

$$\text{Pr}(u | U_l^{(t)}, l \neq k) \approx \frac{1}{M-m+1} \sum_{t=m}^M \text{Pr}(u | U_l^{(t)}, l \neq k)$$

which justify the estimate given by eq (8.5)

(Q3)

To show that bagging estimate $\hat{f}_{\text{bag}}(x)$ converges to the original estimate $f(x)$ as $B \rightarrow \infty$, we need to show that the expected value of $\hat{f}_{\text{bag}}(x)$ converges to $f(x)$ as $B \rightarrow \infty$.

we can write bagging estimate as:

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

where $\hat{f}^{*b}(x)$ is estimate of $f(x)$ from b^{th} bootstrap.

Using Parametric bootstrap.

$$y_i^{*b} = \hat{u}(x_i) + \hat{\epsilon}_i^{*b}$$

where $\hat{u}(x_i)$ is estimate of $u(x)$ from B-spline smoother.

we can then obtain estimate $\hat{f}^{*b}(x)$ by fitting B-Spline smoothers to bootstrap sample y_i^{*b} .

Now, let's Compute expected value $E[\hat{f}_{\text{bag}}(x)]$

$$E[\hat{f}_{\text{bag}}(x)] = E\left[\frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)\right]$$

$$= \frac{1}{B} \sum_{b=1}^B E[\hat{f}^{*b}(x)]$$

we know that B-spline smoothers $\hat{u}(x)$ is a consistent estimator of $u(x)$, meaning that $\hat{u}(x) \rightarrow u(x)$ as $N \rightarrow \infty$ (N is sample size)

$\therefore B \rightarrow \infty$, we can assume that $\hat{u}(x) \rightarrow u(x)$

$$\therefore E[\hat{f}_{\text{bag}}(x)] = \frac{1}{B} \sum_{b=1}^B E[\hat{f}^{*b}(x)]$$

$$\rightarrow \frac{1}{B} \sum_{b=1}^B f(x) = f(x) \text{ as } B \rightarrow \infty$$

(d)

The purpose of the exercise is to investigate the behaviour of a neural network using sigmoid activation function or more generally, any activation function with property of being nearly around the origin, we will take for granted in this exercise that the neural network uses sigmoid activation σ with property that it's almost linear near origin. Let's first recall sigmoid activation σ .

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

Now, let's examine the Taylor expansion for activation σ . $\sigma(v)$ in the context of neural network

$$\begin{aligned}\sigma(x_0^m + x_m^T x) &= \frac{1}{1 + \exp(-x_0^m - x_m^T x)} \\ &= \left(\frac{\exp(x_0^m)}{1 + \exp(x_0^m)} \right) \left(\frac{1}{\exp(-x_m^T x) - 1} \right) \\ &\quad \frac{1}{1 + \exp(x_0^m) + 1}\end{aligned}$$

We can analyze behaviour of activation σ of x approach zero:

$$\exp(-x_m^T x) - 1 \rightarrow -x_m^T x.$$

which leads to:

$$\sigma(x_0^m + x_m^T x) \rightarrow \left(\frac{\exp(x_0^m)}{1 + \exp(x_0^m)} \right) \left(1 + \frac{x_m^T x}{1 + \exp(x_0^m)} \right)$$

From this, it becomes clear that $\sigma(x)$ is essentially a linear of x .

Now, let's define $u(x) = \frac{1}{2}(1+x)$

$$\lim_{x \rightarrow 0} \frac{u(x)}{\sigma(x)} = \lim_{x \rightarrow 0} \frac{1}{2}(1+x)(1+e^{-x})$$

$$= \lim_{x \rightarrow 0} \frac{1}{2}(1+e^{-x}+x+xe^{-x}) = 1$$

considering the case when α_m values are close to zero, the values of $\alpha_m^T x$ will also be small. If we incorporate a bias term and expand π , we obtain $\pi_m = \pi(\alpha_m^T x) \approx -\frac{1}{2}(1 + \alpha_m^T x)$. As a result, the model exhibits a nearly linear behaviour in π since g_k is an identity fn.

In conclusion, we have demonstrated that when employing a sigmoid activation fn. exhibiting near linear properties close to origin, the neural network model tends to approximate a linear behaviour in π . This observation is particularly relevant when the α_m values are close to zero, further emphasizing the linear characteristics of model in such cases.

(Q8)

we will use cross-entropy loss fn. which is defined as

$$L(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

let's derive the forward to backward propagation equations.

forward propagation: this is process of computing the output of neural network given an input x . we calculate output of each layer based on output of previous layers.

initialize $a^{[0]} = x$. for $l = 1, 2, \dots, L$: a. calculate $a^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$. b. calculate $a^{[L]} = g^{[L]}(2^{[L]})$

The final output of network is $\hat{y} = a^{[L]}$

Backward Propagation: This is process of computing the gradients of loss function with respect to weights and biases of neural network. we use the chain rule to compute these gradients.

calculate error in output layer.

$$\begin{aligned} d^{[L]} &= -\frac{\partial L}{\partial z^{[L]}} = \frac{\partial L}{\partial z^{[L]}} = \frac{\partial L}{\partial \hat{y}^{[L]}} \frac{\partial \hat{y}^{[L]}}{\partial z^{[L]}} \\ &= (a^{[L]} - y) \odot g^{[L]}(z^{[L]}) \end{aligned}$$

for $l=L-1, L-2, \dots, 1$: a calculate the error in layer l

$$d^{[l]} = (w^{[l+1]T} f^{[l+1]}) \odot (z^{[l]}) \text{ for } l=1, 2, \dots, L-1$$

Calculate gradients with respect to weights

$$\frac{\partial L}{\partial w^{[l]}} = f^{[l]} a^{[l+1]T} \text{ b. calculate gradients with respect to biases: } \frac{\partial L}{\partial b^{[l]}} = d^{[l]}$$

Once gradient are computed, we can update the weights and biases using a suitable optimization algorithm, such as gradient descent (or any of its variants).

In other words, the context of Cross-entropy loss is defined as

$$R(\theta) = -\sum_{i=1}^N \sum_{k=1}^K y_{ik} \log(f_k(x_i))$$

where y_{ik} is true label of the i th data point for class k , $f_k(x_i)$ is the Predicted Probability for class k given input x_i . The corresponding classifier for this objective ϕ_θ is given by.

$$G(x) = \arg \max_k f_k(x).$$

Let $z_{mij} = \tau(x_{0m} + x_m^T \alpha_i)$ as defined in eq (11.5) in text, let $z_i = (z_{1i} \ z_{2i} \ \dots \ z_{ni})$

Objective of u:

$$R(O) = \sum_{i=1}^N \sum_{k=1}^C (-y_{ik} \log f_{ik}(z_i))$$

$$\frac{\partial R_i}{\partial \beta_{km}} = -\frac{y_{ik}}{f_k(z_i)} g'_k(\beta_k^T z_i) z_{mij}$$

$$\frac{\partial R_i}{\partial \beta_{km}} = -\sum_{k=1}^C \frac{y_{ik}}{f_k(z_i)} g'_k(\beta_k^T z_i) \beta_{km} \sigma'(\alpha_{0m} + \alpha_m^T \alpha_i) z_{mij}$$

Based on these derivates, a gradient descent update can be performed at the (it^{th}) iteration as follows:

$$\begin{aligned} \beta_{km}^{(\text{it}+1)} &= \beta_{km}^{(\text{it})} = -\gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}} \cdot d_{mi}^{(\text{it}+1)} \\ &= \beta_{km}^{(\text{it})} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}} \cdot d_{mi}^{(\text{it})}. \end{aligned}$$

γ_r is learning rate.

By Partial derivative as

$$\frac{\partial R_i}{\partial \beta_{km}} = \text{d}_{mi} \frac{\partial}{\partial \alpha_m} \frac{\partial R_i}{\partial \alpha_m} = S_{mi} z_{mi}$$

From their definitions we have back propagation eq.

$$S_{mi} = \sigma'(\alpha_{0m} + \alpha_m^T \alpha_i) \sum_{k=1}^C \beta_{km} d_{ki}$$

In Summary, forward propagation calculates the Predicted Probabilities for each class given an input, while backward propagation computes the gradient of loss of u. with respect to model Parameters. These gradients are then used to update model Parameters iteratively refining the model's prediction by minimizing cross-entropy loss over training Data