

Assignment 5 by nc41 & cq4

Problem 1 Deep neural networks

Why do deep networks typically outperform shallow networks?

More complex function can be combined in the deep network than in the shallow networks.

As a result, deep networks can be used to find the prediction which might be more similar to actual goal function.

What is leaky ReLU activation and why is it used?

In the normal ReLU, it kills all the nodes whose value is less than 0. It might kill too many nodes which might be useful. Therefore, leaky ReLU gives the nodes, whose value is less than 0, the second chance. These nodes multiply a very small number, such as 0.01, and these nodes can live and transport forward.

In one or more sentences, and using sketches as appropriate, contrast: AlexNet, VGGNet, GoogleNet and ResNet. What is the one defining characteristic of each network?

AlexNet: The Network had a very similar architecture to LeNet, but was deeper, bigger, and featured Convolutional Layers stacked on top of each other (previously it was common to only have a single CONV layer always immediately followed by a POOL layer).

Google Net: Their architecture consisted of a 22 layer deep CNN but reduced the number of parameters from 60 million (AlexNet) to 4 million. Additionally, this paper uses Average Pooling instead of Fully Connected layers at the top of the ConvNet, eliminating a large amount of parameters that do not seem to matter much

VggNet: VGGNet consists of 140 million parameters, which can be a bit challenging to handle.

ResNet: They were able to train a NN with 152 layers while still having lower complexity than VGGNet.

Problem 2 Decision trees, entropy and information gain

2.1

$$H(p) = -p \log p - (1-p) \log (1-p)$$

show $H(p/(p+n))$ $H(s) \leq 1$ $H(s) = 1$ when $p=n$.

When $p=n$

$$H(p/(p+n)) = H\left(\frac{1}{2}\right) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = -\log_2 \frac{1}{2} = \log_2 2 = 1$$

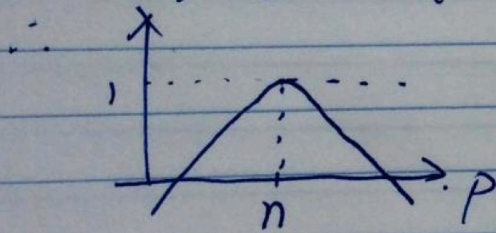
$$H(q) = -q \log q - (1-q) \log (1-q)$$

$$H(q)' = -\log q - 1 + \log (1-q) + 1$$

$$= \log (1-q) - \log (q)$$

when $q < \frac{1}{2}$, $H(q)' > 0$ $q < \frac{1}{2}$ means $p < n$

when $q > \frac{1}{2}$, $H(q)' < 0$ $q > \frac{1}{2}$ means $p > n$



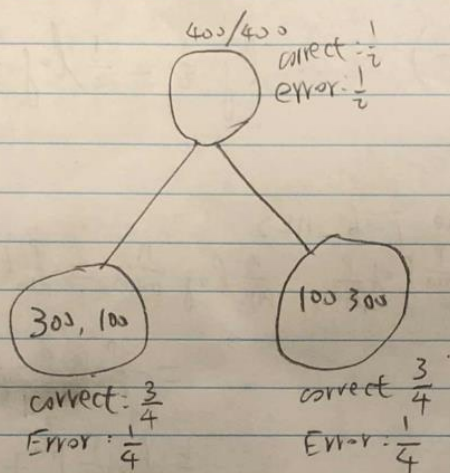
So for any given p and n , $H(s) \leq 1$

2.2 & 2.3

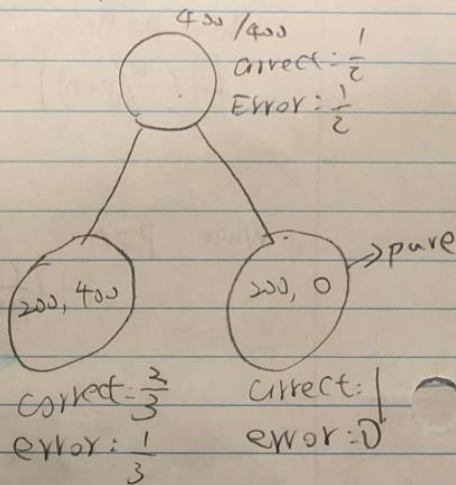
(2)

misclassification rate:

model A



model B



$$\frac{1}{2} - \frac{1}{2} \times \frac{1}{400} \times 100 - \frac{1}{2} \times \frac{1}{400} \times 100$$

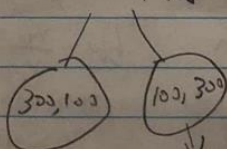
$$= \frac{1}{2} - \frac{1}{4} = \frac{1}{4}$$

$$\frac{1}{2} - \frac{2}{3} \times \frac{1}{600} \times 200 - \frac{1}{2} \times \frac{1}{600} \times 0$$

$$= \frac{1}{2} - \frac{1}{4} = \frac{1}{4}$$

model A and model B have the same information gain:

entropy: model A $H(S) = H(\frac{1}{2}) = 1$ model B the same $1 - H(\frac{1}{2}) = 1$



the same $H(S) = 0.81$

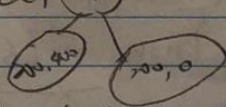
$$-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}$$

$$= -\frac{3}{4} (\log 3 + 2) + \frac{1}{2}$$

$$= 0.81 \rightarrow 1 - \frac{1}{2} \times 0.81 - \frac{1}{2} \times 0.81$$

$$= 1 - 0.405 - 0.405 = 0.19$$

in this case model B > model A



$$-\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$$

$$= -\frac{2}{3} (\log 2 - \log 3) + \frac{1}{3} \log 3$$

$$= \log 3 - \frac{2}{3} = 1.58 - 0.67 = 0.91$$

$$= 0.91$$

$$1 - \frac{3}{4} \times 0.91 - 0 = 0.32$$

Gini Index:

Model A $\Rightarrow 2 \times \frac{1}{2} (1 - \frac{1}{2})$
 $= \frac{1}{2}$

Diagram for Model A: A root node splits into two nodes, (300, 100) and (100, 300). Each node has a weight of $\frac{3}{4}$ below it. Below (300, 100) is the calculation $2 \times \frac{3}{4} (1 - \frac{3}{4}) = \frac{3}{8}$. Below (100, 300) is the calculation $2 \times \frac{3}{4} (1 - \frac{3}{4}) = \frac{3}{8}$. At the bottom is the calculation $1 - \frac{1}{2} \times \frac{3}{8} - \frac{1}{2} \times \frac{3}{8} = \frac{5}{8} = 0.625$.

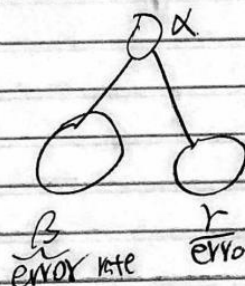
Model B $\Rightarrow 2 \times \frac{1}{2} (1 - \frac{1}{2})$

Diagram for Model B: A root node splits into two nodes, (200, 400) and (200, 0). The left node has a weight of $\frac{2}{3}$ below it, and the right node has a weight of $\frac{1}{3}$ below it. Below (200, 400) is the calculation $2 \times \frac{2}{3} \times \frac{1}{3} = \frac{4}{9}$. Below (200, 0) is the calculation 0 . At the bottom is the calculation $1 - \frac{2}{3} \times \frac{4}{9} - 0 = \frac{2}{3} = 0.667$.

In Gini Index information gain Model B > Model A

So we choose Model B

(3) If root misclassification rate is α



When splitting, the information gain > 0 means we can have progress when splitting

because $\beta, \gamma \leq \frac{1}{2}$
 $\text{error rate} = 1 - \max(p, 1-p)$

$$\alpha - \frac{|D_L|}{|D|} \cdot \beta - \frac{|D_R|}{|D|} \cdot \gamma$$

Assume $\Rightarrow \alpha = \frac{1}{100} \beta + \frac{99}{100} \gamma \Rightarrow$ the information gain still can be more than 0

if $\beta > \alpha$ $\gamma < \alpha$

So we still split, although β error rate will be more than α if the data in that branch is not larger enough to make information gain less than 0 the tree will still split.

Problem 3 Bagging

hw 5 bagging

$$(1) \quad E_x[\{f(x) + \epsilon_1(x) - f(x)\}^2] = E_x[\epsilon_1(x)^2]$$

$$(h_1(x) = f(x) + \epsilon_1(x))$$

$$E_{av} = \frac{1}{L} \sum_{i=1}^L E_x[\epsilon_i(x)^2]$$

$$h_{bag}(x) = \frac{1}{L} \sum_{i=1}^L h_i(x) \rightarrow \text{vote}$$

$$\begin{aligned} \epsilon_{bag}(x) &= h_{bag}(x) - f(x) = \frac{1}{L} \sum_{i=1}^L h_i(x) - f(x) = \frac{1}{L} \sum_{i=1}^L (f(x) + \epsilon_i(x)) - f(x) \\ &= \frac{1}{L} \sum_{i=1}^L \epsilon_i(x) \end{aligned}$$

$$E_{bag} = E[\epsilon_{bag}(x)^2]$$

$$\text{if } E_x[\epsilon_i(x)] = 0 \quad E_x[\epsilon_m(x)\epsilon_n(x)] = 0 \quad \text{when } m \neq n$$

$$\begin{aligned} E_{bag} &= E[\epsilon_{bag}(x)^2] = E\left[\frac{1}{L} \left(\sum_{i=1}^L \epsilon_i(x)\right)^2\right] \\ &= \frac{1}{L^2} E\left[\sum_{i=1}^L \epsilon_i(x)^2 + \sum_{i \neq m} \epsilon_i(x)\epsilon_m(x)\right] \\ &= \frac{1}{L^2} E\left[\sum_{i=1}^L \epsilon_i(x)^2\right] \\ &= \frac{1}{L} \cdot \frac{1}{L} \sum_{i=1}^L E[\epsilon_i(x)^2] \\ &= \frac{1}{L} \cdot E_{av} \end{aligned}$$

$$\begin{aligned}
 f(x) &= x^2 \text{ here, } \lambda_i = \frac{1}{L} \\
 f\left(\sum_i \frac{1}{L} \epsilon_i\right) &\leq \sum_i \frac{1}{L} f(\epsilon_i) \\
 \Rightarrow \left(\sum_i \frac{1}{L} \epsilon_i\right)^2 &\leq \sum_i \frac{1}{L} \epsilon_i^2 \\
 \Rightarrow E_{\text{bag}} &\leq E_{\text{av}}
 \end{aligned}$$

Extra

In this section, we tried several different architectures, firstly using a small 6 layer net. The performance is bad. And since training on CPU is really slow. It is difficult to tune hyperparameters. We then implement it using Keras, and ran it on colab. However, later Devika said we can not use framework but just the API. We have to re-implement it using API.

So below is the result we achieved. And you can see more details from the two htmls file.

	Train Accuracy	Validation Accuracy	Test Accuracy
Keras(extra.html)	71.17%	73.4%	73.7%
Api(Conv....html)	94.3%	74.6 %	72.9 %