

Problem Set 5

Zhuo Li(zhuol2)

Handed In: April 10, 2017

1. Neural Networks

- (a) For each training example d , we know that we need to update weights on the opposite direction of the gradient (gradient descent). So we have

$$\Delta\omega_{ij} = -R \frac{\partial E_d}{\partial \omega_{ij}} \quad (1)$$

ω_{ij} can only influence the output through net_j , such that

$$\text{net}_j = \sum \omega_{ij} x_{ij} \quad (2)$$

where x_{ij} is from the previous layer of unit j . From (1)(2) we have

$$\frac{\partial E_d}{\partial \omega_{ij}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial \omega_{ij}} = \frac{\partial E_d}{\partial \text{net}_j} x_{ij} \quad (3)$$

$$\Delta\omega_{ij} = -R \frac{\partial E_d}{\partial \text{net}_j} x_{ij} \quad (4)$$

Then we consider two cases:

- i. **Unit j is the output unit.** So net_j can only influence the rest of the network through o_j . So

$$\frac{\partial E_d}{\partial \text{net}_j} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \quad (5)$$

We also have

$$E_d = \frac{1}{2} \sum_{k \in K} (t_k - o_k)^2 \quad (6)$$

$$o_j = \max(0, \text{net}_j) \quad (7)$$

Thus

$$\frac{\partial o_j}{\partial \text{net}_j} = \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ 1 & \text{otherwise} \end{cases} \quad (8)$$

Then

$$\begin{aligned} \frac{\partial E_d}{\partial \text{net}_j} &= \frac{1}{2} 2(t_j - o_j) \frac{\partial (t_j - o_j)}{\partial \text{net}_j} \\ &= \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ -o_j(t_j - o_j) & \text{otherwise} \end{cases} \end{aligned} \quad (9)$$

So

$$\begin{aligned}\Delta\omega_{ij} &= -R \frac{\partial E_d}{\partial \text{net}_j} x_{ij} = R\delta_j x_{ij} \\ &= \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ R(t_j - o_j)o_j x_{ij} & \text{otherwise} \end{cases} \end{aligned} \quad (10)$$

where

$$\delta_j \equiv \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ (t_j - o_j)o_j & \text{otherwise} \end{cases} \quad (11)$$

ii. **Unit j is the hidden unit.** So net_j can influence the network only through $\text{downstream}(j)$. Then

$$\begin{aligned}\frac{\partial E_d}{\partial \text{net}_j} &= \sum_{k \in \text{downstream}(j)} \frac{\partial E_d}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial \text{net}_j} \\ &= \sum_{k \in \text{downstream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial \text{net}_j} \\ &= \sum_{k \in \text{downstream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \\ &= \sum_{k \in \text{downstream}(j)} -\delta_k \omega_{jk} \frac{\partial o_j}{\partial \text{net}_j} \end{aligned} \quad (12)$$

Remember what we have in (8)

$$\frac{\partial E_d}{\partial \text{net}_j} = \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ \sum_{k \in \text{downstream}(j)} -\delta_k \omega_{jk} & \text{otherwise} \end{cases} \quad (13)$$

So

$$\begin{aligned}\Delta\omega_{ij} &= -R \frac{\partial E_d}{\partial \text{net}_j} x_{ij} = R\delta_j x_{ij} \\ &= \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ Rx_{ij} \sum_{k \in \text{downstream}(j)} \delta_k \omega_{jk} & \text{otherwise} \end{cases} \end{aligned} \quad (14)$$

where

$$\delta_j \equiv \begin{cases} 0 & \text{if } \text{net}_j \leq 0 \\ \sum_{k \in \text{downstream}(j)} \delta_k \omega_{jk} & \text{otherwise} \end{cases} \quad (15)$$

(b) i. From (6) we have the gradient of squared loss with respect to output o_k is

$$\frac{\partial E_d}{\partial o_k} = (o_k - t_k)$$

So we can have the following codes:

```

1 def squared_loss_gradient (output_activations , y):
2     #IMPLEMENT THIS
3     (m, n) = output_activations.shape
4     gradient = np.zeros((m, n))
5     for i in range(m):
6         for j in range(n):
7             gradient[i, j] = output_activations[i, j] -
                        y[i, j]
8     return gradient
9 #endDef

```

And (8) is the formula for the derivative of the rectifier function, then

```

1 def relu_derivative (z):
2     #IMPLEMENT THIS!
3     (m, n) = z.shape
4     derivative = np.zeros((m, n))
5     for i in range(m):
6         for j in range(n):
7             if z[i, j] > 0:
8                 derivative[i, j] = 1
9             else:
10                 derivative[i, j] = 0
11     return derivative
12 #endDef

```

ii. We can run the *parameter_tuning.py*, and get the following results for **circles**:

```

1 ===== circles =====
2 — size = 10 , R = 0.1 , func = relu , node = 10 —
3 <Average Accuracy> 100.0
4 — size = 10 , R = 0.1 , func = relu , node = 50 —
5 <Average Accuracy> 100.0
6 — size = 10 , R = 0.1 , func = tanh , node = 10 —
7 <Average Accuracy> 100.0
8 — size = 10 , R = 0.1 , func = tanh , node = 50 —
9 <Average Accuracy> 100.0
10 — size = 10 , R = 0.01 , func = relu , node = 10 —
11 <Average Accuracy> 100.0
12 — size = 10 , R = 0.01 , func = relu , node = 50 —
13 <Average Accuracy> 100.0
14 — size = 10 , R = 0.01 , func = tanh , node = 10 —
15 <Average Accuracy> 52.7647714604
16 — size = 10 , R = 0.01 , func = tanh , node = 50 —
17 <Average Accuracy> 50.3448001274
18 — size = 50 , R = 0.1 , func = relu , node = 10 —

```

```

19 <Average Accuracy> 100.0
20 — size = 50 , R = 0.1 , func = relu , node = 50 —
21 <Average Accuracy> 100.0
22 — size = 50 , R = 0.1 , func = tanh , node = 10 —
23 <Average Accuracy> 66.0081223125
24 — size = 50 , R = 0.1 , func = tanh , node = 50 —
25 <Average Accuracy> 57.1444497531
26 — size = 50 , R = 0.01 , func = relu , node = 10 —
27 <Average Accuracy> 74.1025641026
28 — size = 50 , R = 0.01 , func = relu , node = 50 —
29 <Average Accuracy> 99.6273291925
30 — size = 50 , R = 0.01 , func = tanh , node = 10 —
31 <Average Accuracy> 51.9190953974
32 — size = 50 , R = 0.01 , func = tanh , node = 50 —
33 <Average Accuracy> 51.2581621277
34 — size = 100 , R = 0.1 , func = relu , node = 10 —
35 <Average Accuracy> 100.0
36 — size = 100 , R = 0.1 , func = relu , node = 50 —
37 <Average Accuracy> 100.0
38 — size = 100 , R = 0.1 , func = tanh , node = 10 —
39 <Average Accuracy> 52.0074852684
40 — size = 100 , R = 0.1 , func = tanh , node = 50 —
41 <Average Accuracy> 51.8633540373
42 — size = 100 , R = 0.01 , func = relu , node = 10 —
43 <Average Accuracy> 65.0947603122
44 — size = 100 , R = 0.01 , func = relu , node = 50 —
45 <Average Accuracy> 74.7356266921
46 — size = 100 , R = 0.01 , func = tanh , node = 10 —
47 <Average Accuracy> 52.668418538
48 — size = 100 , R = 0.01 , func = tanh , node = 50 —
49 <Average Accuracy> 51.2501990763

```

Then we have the optimal parameters for **circles**:

```

[Optimal Batch Size] 10
[Optimal Learning Rate] 0.1
[Optimal Activation Function] relu
[Optimal Hidden Layer Width] 10
[Maximum Accuracy] 100.0%

```

And results for **mnist**:

```

1 ===== mnist =====
2 — size = 10 , R = 0.1 , func = relu , node = 10 —
3 <Average Accuracy> 95.9688541337
4 — size = 10 , R = 0.1 , func = relu , node = 50 —
5 <Average Accuracy> 87.5582991689
6 — size = 10 , R = 0.1 , func = tanh , node = 10 —

```

```

7 <Average Accuracy> 96.9621801688
8 — size = 10 , R = 0.1 , func = tanh , node = 50 —
9 <Average Accuracy> 96.9038261068
10 — size = 10 , R = 0.01 , func = relu , node = 10 —
11 <Average Accuracy> 96.6199395025
12 — size = 10 , R = 0.01 , func = relu , node = 50 —
13 <Average Accuracy> 96.4864811173
14 — size = 10 , R = 0.01 , func = tanh , node = 10 —
15 <Average Accuracy> 96.8787634501
16 — size = 10 , R = 0.01 , func = tanh , node = 50 —
17 <Average Accuracy> 96.1944075876
18 — size = 50 , R = 0.1 , func = relu , node = 10 —
19 <Average Accuracy> 96.2193238624
20 — size = 50 , R = 0.1 , func = relu , node = 50 —
21 <Average Accuracy> 96.4613139022
22 — size = 50 , R = 0.1 , func = tanh , node = 10 —
23 <Average Accuracy> 96.9455135544
24 — size = 50 , R = 0.1 , func = tanh , node = 50 —
25 <Average Accuracy> 96.5115019505
26 — size = 50 , R = 0.01 , func = relu , node = 10 —
27 <Average Accuracy> 96.6199604142
28 — size = 50 , R = 0.01 , func = relu , node = 50 —
29 <Average Accuracy> 96.6700439041
30 — size = 50 , R = 0.01 , func = tanh , node = 10 —
31 <Average Accuracy> 96.394689268
32 — size = 50 , R = 0.01 , func = tanh , node = 50 —
33 <Average Accuracy> 96.261084501
34 — size = 100 , R = 0.1 , func = relu , node = 10 —
35 <Average Accuracy> 96.3945847096
36 — size = 100 , R = 0.1 , func = relu , node = 50 —
37 <Average Accuracy> 96.5949291251
38 — size = 100 , R = 0.1 , func = tanh , node = 10 —
39 <Average Accuracy> 96.8787320826
40 — size = 100 , R = 0.1 , func = tanh , node = 50 —
41 <Average Accuracy> 96.2945431999
42 — size = 100 , R = 0.01 , func = relu , node = 10 —
43 <Average Accuracy> 96.7201378498
44 — size = 100 , R = 0.01 , func = relu , node = 50 —
45 <Average Accuracy> 96.761856665
46 — size = 100 , R = 0.01 , func = tanh , node = 10 —
47 <Average Accuracy> 96.2527930172
48 — size = 100 , R = 0.01 , func = tanh , node = 50 —
49 <Average Accuracy> 96.2193865975

```

Then we have the optimal parameters for **circles**:

[Optimal Batch Size] 10
 [Optimal Learning Rate] 0.1
 [Optimal Activation Function] tanh
 [Optimal Hidden Layer Width] 10
 [Maximum Accuracy] 96.9621801688%

- iii. Run the *plot_learning_curve.py*, we can get learning curves in Figure 1 and Figure 2.

Figure 1: Learning Curve for Circles

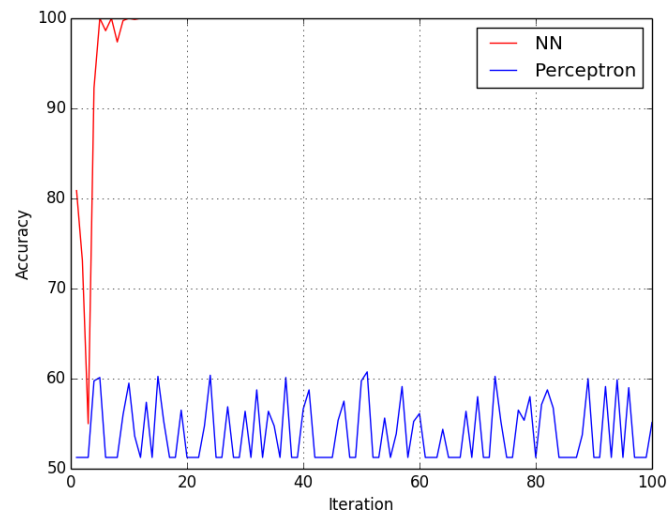
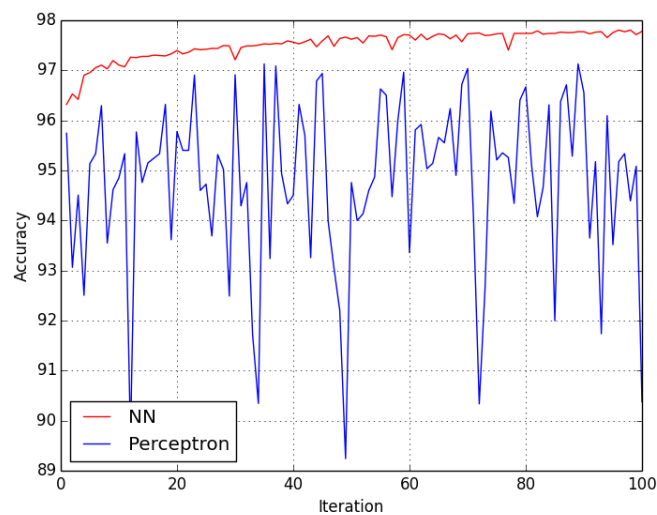


Figure 2: Learning Curve for MNIST



- In Figure 1, we can find that the performance of Perceptron is terrible with accuracy between 50% and 60%. This is because the circle dataset is

not linearly separable, while Perceptron can only learn linear functions. Neural network performs much better and reach nearly 100% accuracy after a few iterations.

- In Figure 2, Perceptron can learn quite well with the average accuracy around 95%, but it is still not as good as neural network. In this case, we can also find that it takes more iterations for neural networks converging.

2. Multi-class classification

- (a)
 - i.
 - **OvA**: We learn k classifiers.
 - **AvA**: We learn $\binom{k}{2}$ classifiers.
 - ii.
 - **OvA**: We use m examples to learn each classifier.
 - **AvA**: We use $\frac{2m}{k}$ examples to learn each classifier.
 - iii.
 - **OvA**: We use the winner takes all (WTA) strategy, such that $f(x) = \operatorname{argmax}_i \omega_i^T x$. The "score" $\omega_i^T x$ can be thought of as the probability that x has label i .
 - **AvA**: We have two options to make a decision. One is to classify example x to take label i if i wins on x more often than any $j = 1, \dots, k$. Alternatively, we can do a tournament. Starting with $n/2$ pairs, continue with the winners and go down iteratively.
 - iv.
 - **OvA**: We need to train $O(k)$ (linear) number of classifiers.
 - **AvA**: We need to train $O(k^2)$ (quadratic) number of classifiers.
- (b) I prefer **OvA** than **AvA**. Because we have better computational complexity of **OvA** over **AvA**. **OvA** has more examples to learn than **AvA** as well.
- (c) The analysis above doesn't hold for **Kernel Perceptron**. In this case, we need dual representations, **AvA** has smaller learning problems in terms of number of examples, thus being preferable when ran in dual. And it is also more expressive. So I prefer **AvA** in this case.
- (d)
 - **OvA**: We have m examples supplied for each classifier and k classifiers to learn. So the overall complexity is

$$O(dm^2)k = O(kdm^2)$$

- **AvA**: We have $\frac{m}{k}$ examples supplied for each classifier and $\binom{k}{2}$ classifiers to learn. So the overall complexity is

$$O(d(\frac{m}{k})^2) \binom{k}{2} = O(\frac{dm^2}{k^2})O(k^2) = O(dm^2)$$

So **AvA** is the most efficient.

- (e)
 - **OvA**: We have m examples supplied for each classifier and k classifiers to learn. So the overall complexity is

$$O(d^2m)k = O(kd^2m)$$

- **AvA**: We have $\frac{m}{k}$ examples supplied for each classifier and $\binom{k}{2}$ classifiers to learn. So the overall complexity is

$$O(d^2(\frac{m}{k}))\binom{k}{2} = O(\frac{d^2m}{k})O(k^2) = O(kd^2m)$$

So **OvA** and **AvA** are in the same order of time complexity. We cannot say which one is the most efficient.

- (f) • **Counting**: For each example, we need first predict with all the classifiers, which take $m(m-1)/2$, i.e. $O(m^2)$, time complexity. Then we do a majority vote, which takes $O(1)$ time. So the overall time complexity is $O(m^2)$.
- **Knockout**: For each round, we knocked out half of the classifiers. So the overall time complexity is $O(\log m)$.

3. Probability Review

- (a) i. Let X be the random variable donating the number of children in a family.
- Town A: Since each family has just one child. So

$$E[X] = 1$$

- Town B: We have

$$\begin{aligned} E[X] &= \sum_{k=1}^{\infty} k \times P(X = k) \\ &= 1 \times \frac{1}{2} + 2 \times \left(\frac{1}{2}\right)^2 + 3 \times \left(\frac{1}{2}\right)^3 + \dots \end{aligned} \tag{1}$$

To compute this, we multiply both sides of (1) by 2.

$$2E[X] = 1 + 2 \times \left(\frac{1}{2}\right) + 3 \times \left(\frac{1}{2}\right)^2 + \dots \tag{2}$$

Subtract both sides of (2) with corresponding sides of (1).

$$\begin{aligned} E[X] &= 1 + \frac{1}{2} + \left(\frac{1}{2}\right)^2 + \dots \\ &= \lim_{k \rightarrow \infty} \frac{1 \times (1 - (\frac{1}{2})^k)}{1 - \frac{1}{2}} \\ &= 2 \end{aligned} \tag{3}$$

- ii. Let X, Y be the random variable donating the numbers of boys and girls, respectively, in the town at the end of one generation. And there are N_A families in Town A and N_B in Town B.

- Town A: Since each family has just one child with equal possibility to be a boy or a girl. So

$$E[X] = E[Y] = N_A/2$$

Then

$$\frac{E[X]}{E[Y]} = \frac{N_A/2}{N_A/2} = 1$$

- Town B: Each family will have exactly one boy. So

$$E[X] = N_B$$

For girls

$$\begin{aligned} E[Y] &= N_B \\ &= N_B \left(\sum_{k=0}^{\infty} k \times P(Z = k) \right) \\ &= N_B \left(1 \times \left(\frac{1}{2}\right)^2 + 2 \times \left(\frac{1}{2}\right)^3 + 3 \times \left(\frac{1}{2}\right)^4 + \dots \right) \end{aligned} \quad (1)$$

where Z is the random variable donating number of girls in each family. Use the same method as in (i)

$$2E[Y] = N_B \left(1 \times \frac{1}{2} + 2 \times \left(\frac{1}{2}\right)^2 + 3 \times \left(\frac{1}{2}\right)^3 + \dots \right) \quad (2)$$

From (1)(2), we have

$$\begin{aligned} E[Y] &= N_B \left(\frac{1}{2} + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^3 + \dots \right) \\ &= N_B \lim_{k \rightarrow \infty} \frac{\frac{1}{2} \left(1 - \left(\frac{1}{2}\right)^k \right)}{1 - \frac{1}{2}} \\ &= N_B \end{aligned} \quad (3)$$

Then

$$\frac{E[X]}{E[Y]} = \frac{N_B}{N_B} = 1$$

- (b) i. Proof: From the chain rule, we have

$$P(A, B) = P(A|B)P(B) \quad (1)$$

$$P(A, B) = P(B|A)P(A) \quad (2)$$

So from (1)(2) we will have

$$P(A|B)P(B) = P(B|A)P(A)$$

Then

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

ii. From the chain rule, we can rewrite as

$$P(A, B, C) = P(A|B, C)P(B, C) = P(A|B, C)P(B|C)P(C)$$

(c) Proof: According to the definition of expectation:

$$\begin{aligned} E[X] &= 1 \times P(X = 1) + 0 \times P(X = 0) \\ &= 1 \times P(A) \\ &= P(A) \end{aligned}$$

(d) i. No.

We can calculate

$$P(X = 0) = 1/15 + 1/10 + 4/15 + 8/45 = 11/18$$

$$P(Y = 0) = 1/15 + 1/15 + 4/15 + 2/15 = 8/15$$

And

$$P(X = 0, Y = 0) = 1/15 + 4/15 = 1/3$$

Then

$$P(X = 0, Y = 0) \neq P(X = 0)P(Y = 0)$$

So X is not independent of Y .

ii. Yes. We need to prove for any $x, y, z \in \{0, 1\}$:

$$P(X = x, Y = y|Z = z) = P(X = x|Z = z)P(Y = y|Z = z)$$

We have:

$$P(Z = 0) = 1/15 + 1/15 + 1/10 + 1/0 = 1/3 \quad P(Z = 1) = 1 - P(Z = 0) = 2/3 \quad (1)$$

We have:

$$P(X = 0|Z = 0) = \frac{1/15 + 1/10}{1/3} = 1/2 \quad (2)$$

$$P(X = 1|Z = 0) = \frac{1/15 + 1/10}{1/3} = 1/2 \quad (3)$$

$$P(Y = 0|Z = 0) = \frac{1/15 + 1/15}{1/3} = 2/5 \quad (4)$$

$$P(Y = 1|Z = 0) = \frac{1/10 + 1/10}{1/3} = 3/5 \quad (5)$$

$$P(X = 0|Z = 1) = \frac{4/15 + 8/45}{2/3} = 2/3 \quad (6)$$

$$P(X = 1|Z = 1) = \frac{2/15 + 4/45}{2/3} = 1/3 \quad (7)$$

$$P(Y = 0|Z = 1) = \frac{4/15 + 2/15}{2/3} = 3/5 \quad (8)$$

$$P(Y = 1|Z = 1) = \frac{8/45 + 4/45}{2/3} = 2/5 \quad (9)$$

And:

$$\begin{aligned} P(X = 0, Y = 0|Z = 0) &= \frac{1/15}{1/3} = 1/5 \\ &= P(X = 0|Z = 0)P(Y = 0|Z = 0) \end{aligned} \quad (10)$$

$$\begin{aligned} P(X = 0, Y = 1|Z = 0) &= \frac{1/10}{1/3} = 3/10 \\ &= P(X = 0|Z = 0)P(Y = 1|Z = 0) \end{aligned} \quad (11)$$

$$\begin{aligned} P(X = 1, Y = 0|Z = 0) &= \frac{1/15}{1/3} = 1/5 \\ &= P(X = 1|Z = 0)P(Y = 0|Z = 0) \end{aligned} \quad (12)$$

$$\begin{aligned} P(X = 1, Y = 1|Z = 0) &= \frac{1/10}{1/3} = 3/10 \\ &= P(X = 1|Z = 0)P(Y = 1|Z = 0) \end{aligned} \quad (13)$$

$$\begin{aligned} P(X = 0, Y = 0|Z = 1) &= \frac{4/15}{2/3} = 2/5 \\ &= P(X = 0|Z = 1)P(Y = 0|Z = 1) \end{aligned} \quad (14)$$

$$\begin{aligned} P(X = 0, Y = 1|Z = 1) &= \frac{8/45}{2/3} = 4/15 \\ &= P(X = 0|Z = 1)P(Y = 1|Z = 1) \end{aligned} \quad (15)$$

$$\begin{aligned} P(X = 1, Y = 0|Z = 1) &= \frac{2/15}{2/3} = 1/5 \\ &= P(X = 1|Z = 1)P(Y = 0|Z = 1) \end{aligned} \quad (16)$$

$$\begin{aligned} P(X = 1, Y = 1|Z = 1) &= \frac{4/45}{2/3} = 2/15 \\ &= P(X = 1|Z = 1)P(Y = 1|Z = 1) \end{aligned} \quad (17)$$

So X is conditionally independent of Y given Z .

iii.

$$\begin{aligned} P(X = 0|X + Y > 0) &= \frac{P(X = 0, X + Y > 0)}{P(X + Y > 0)} \\ &= \frac{1/10 + 8/45}{1 - 1/15 - 4/15} \\ &= 5/12 \end{aligned}$$