

ESL: 11.1
 From The Neural Network Side, Since there is ~~not~~ one hidden layer.

$$f_k(x) = \beta_0 + \beta_k^T \theta(\alpha_{0k} + \alpha_k^T x)$$

From The PPR side we know:

$$f(x) = \sum_{n=1}^N g_n(w_n^T x)$$

Also we know:

$$g_n(w_n^T x) = \beta_n \theta(\alpha_{0n} + s_n(w_n^T x))$$

Then since is a PPR, we know that $w = \frac{\hat{w}}{\|\hat{w}\|}$,

Then, we need to find some s_n and \hat{w} , such that will allow us to equate both weights; this is $s = \|\alpha_n\|$ and

$$w_n^T = \frac{\alpha_n}{\|\alpha_n\|}$$

Then $g_n(w_n^T x) = f_n(x) = \beta_0 + \beta_n \theta(\alpha_{0n} + \alpha_n^T x)$ And

$$f(x) = \sum_{n=1}^K f_n(x)$$

The same argument could be done if we use a logistic regression for the $g_n(w_n^T x)$ of PPR and a classification network from the neural network. Where $g(t) = \frac{e^t}{\sum e^t}$.

ESL: 11.2

If the function is linear we know that $f'' = 0$.

Then if $y = \frac{1}{1+e^{-v}}$

$$\frac{\partial y}{\partial v} = \frac{e^{-v}}{(1+e^{-v})^2}$$

$$\frac{\partial^2 y}{\partial v^2} = e^{-v}(1+e^{-v})^2(2 \cdot (1+e^{-v})^{-1} - 1)$$

Then if $v=0$.

$$\left. \frac{\partial^2 y}{\partial v^2} \right|_{v=0} = \frac{1}{4} \cdot \left(\frac{2}{2} - 1 \right)$$

$= 0 // \Rightarrow$ Then we could argue that
is linear at $v=0$.

ESL 11.4

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

$$f_k(x) = \frac{e^{T_k}}{\sum_{k=1}^K e^{T_k}}$$

Since there is no hidden layer

$$T_k = \beta_k + \beta_k^T X$$

AND NOW WE ASSUME THAT THERE $K-1$ IS BASE CATEGORIES. THEN

$$T_k - T_1 = (\beta_k - \beta_1) + (\beta_k - \beta_1)^T X$$

$$\therefore \Pr(K=k) = f_k(x) = \frac{e^{T_k - T_1}}{1 + e^{T_k - T_1}}$$

With transform for multiple K AND USING CROSS ENTROPY.

$$- \sum y_{ik} \log(f_k(x))$$

\Rightarrow WITH IS THE MULTINOmIAL OBJECTIVE FUNCTION.

S.1 Bishop

First we need to find the conversion or relation between $\tanh(a)$ and $\theta(a)$.

$$\frac{e^a - e^{-a}}{e^a + e^{-a}} = \frac{e^a}{e^a + e^{-a}} - \frac{e^{-a}}{e^a + e^{-a}} = \frac{e^a}{e^a + e^{-a}} - \frac{e^{-a} + e^a - e^a}{e^a + e^{-a}}$$
$$= \frac{e^a}{e^a + e^{-a}} - \left(1 - \frac{e^a}{e^a + e^{-a}} \right)$$

$$= \frac{2e^a}{e^a + e^{-a}} - 1$$

$$= \frac{2}{1 + e^{-2a}} - 1$$

$$\frac{e^a - e^{-a}}{e^a + e^{-a}}$$

$$= 2 \theta(2a) - 1$$

Then if we replace this into

$$a_k = \sum_{j=1}^n w_{kj}^{(1)} \cdot (2\theta(2a) - 1) + w_{k0}^{(1)}$$

$$a_k = \sum_{j=1}^n 2w_{kj}^{(1)} \theta(2a) - 2w_{kj}^{(1)} + w_{k0}^{(1)}$$

$$a_k = \sum_j 2w_{kj}^{(1)} \theta(2a) + \sum_j [-2w_{kj}^{(1)} + w_{k0}^{(1)}]$$

Then if for a regular neural network: ~~we have to~~ we have to convert into:

$$\omega_p^{(2p)} = 2\omega_R^{(2)}, \quad a_j^{(1p)} = 2a_j^{(2)}, \quad \omega_{k0}^{(2p)} = -\sum_{j=1}^n \omega_{kj}^{(2)} + \omega_{k0}^2$$

Using the eqn. 1.2, we transform a_R into

$$a_p = \sum_{j=1}^n \omega_p^{(2p)} \theta(a) + \omega_{p0}^{(2p)}$$

Bishop 5.5

$$P(\tau/w_1, \dots, w_K) = \prod_{k=1}^K y_{\tau k}^{\tau_k} \Rightarrow \text{network output for } K \text{ class.}$$

Then for N DATA points we have that

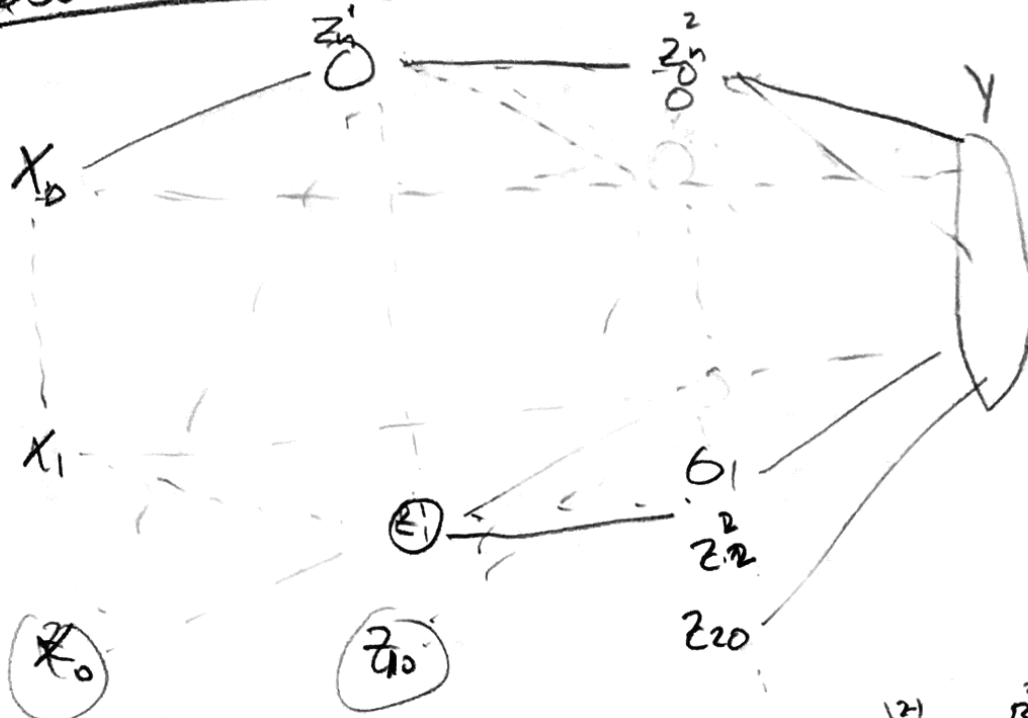
$$P(T/w_1, \dots, w_K) = \prod_{k=1}^K \prod_{n=1}^N y_{nk}^{\tau_{nk}}$$

Then if we compare for the log function.

$$l(T/w_1, \dots, w_K) = - \sum_{k=1}^K \sum_{n=1}^N \tau_{nk} \cdot \log(y_{nk}).$$

Which is the log-likelihood function of the
5.24.

Question 6



Equation:

$$z_n^1 = \sigma(\alpha_{0n} + \alpha_n^T X)$$

$$\hat{y}_k = \sigma(\beta_{0k} + \beta_k^T z^1 + \delta_k^T X + \delta_{0k})$$

$$R = \frac{1}{2} \sum_{k=1}^K (\hat{y}_k - y)^2$$

$$z_n^2 = \sigma(\beta_{0n}^{(2)} + \beta^{(2)T} z_n^1)$$

$$\frac{\partial R}{\partial \beta_k} = \delta_{k1} z_{n1}^1, \dots, \quad s_{n1} = \sigma'(\alpha_n^T X_n) \sum_{k=1}^K \beta_{kn} \delta_{kn}$$

δ_{kn} are errors in output layer

$$\frac{\partial R}{\partial \alpha_n} = s_{n1} x_{i1}^1, \dots$$

$$\frac{\partial R}{\partial \delta_k} = 2(y_n - \hat{y}_n) \cdot \hat{y}_n (1 - \hat{y}_n) \cdot X$$

* The main change is that we have another set of parameters that go from the input layer to the output layer, creating a new partial derivative.

hw8 q7

June 1, 2020

1 Question 7

1.1 Part A

```
[60]: import numpy as np

lambda_list = [0.001,0.01,0.1,.25,.5,.75,1,2,3,5]
hiddenSize = 10
epoch_list = list(range(0,10000,100))
hiddenSize_list = list(range(1,11,1))

def sigmoid(x):
    rv = 1/(1+np.exp(-x))
    return rv

def deriv_sigmoid(output):
    return output*(1-output)

def predict(X,alpha,beta):
    l_0 = X
    l_1 = sigmoid(np.dot(l_0,alpha))
    l_2 = sigmoid(np.dot(l_1,beta))
    return l_2

def Neural_backprop(itera, hiddenSize, X, y, eta, Lambda):

    alpha_t = 2*np.random.random((3,hiddenSize)) - 1
    beta_t = 2*np.random.random((hiddenSize,1)) - 1
    alpha_t_1 = alpha_t.copy()
    beta_t_1 = beta_t.copy()
    for j in range(itera):
        l_0 = X
        l_1 = sigmoid(np.dot(l_0,alpha_t))
        l_2 = sigmoid(np.dot(l_1,beta_t))
        l_2_error = l_2 - y
```



```

l_2_delta = l_2_error*deriv_sigmoid(l_2)
l_1_error = l_2_delta.dot(beta_t.T)
l_1_delta = l_1_error * deriv_sigmoid(l_1)
beta_t= beta_t_1*(1- 2*Lambda*eta) + eta * (l_1.T.dot(l_2_delta))
alpha_t = alpha_t_1*(1- 2*Lambda*eta) + eta*(l_0.T.dot(l_1_delta))
#mean_sq_error = (l_2_error.copy())**2
#mean_sq_error = np.mean(mean_sq_error)
return alpha_t,beta_t

```

[]:

1.2 Part B

```

[8]: def dgp(n):
    x1 = np.random.normal(0,1,n)
    x2 = np.random.normal(0,1,n)
    X = np.concatenate((x1.reshape(n,1),x2.reshape(n,1)),axis=1)
    Z = np.random.normal(0,1,n)
    a2 = np.zeros((2,1))
    a2[0][0]=3
    a2[1][0]=-3
    Y=sigmoid(3*X@np.ones((2,1))) +(X@a2)**2 + .3*Z.reshape(n,1)
    return (Y,X)
y_train, X_train = dgp(100)
X_train= np.concatenate((X_train,np.ones((100,1))),axis=1)
y_test, X_test = dgp(1000)
X_test= np.concatenate((X_test,np.ones((1000,1))),axis=1)

```

```

[57]: error_train_6000 = []
error_test_6000 = []
for i in lambda_list:
    alpha, beta = Neural_backprop(6000, 10, X_train, y_train, 0.01,i)
    sq_error_test_mean = np.mean((y_test - predict(X_test,alpha,beta))**2)
    sq_error_train_mean = np.mean((y_train - predict(X_train,alpha,beta))**2)
    error_test_6000.append(sq_error_test_mean)
    error_train_6000.append(sq_error_train_mean)

error_train_4000 = []
error_test_4000 = []
for i in lambda_list:
    alpha, beta = Neural_backprop(4000, 10, X_train, y_train, 0.01,i)
    sq_error_test_mean = np.mean((y_test - predict(X_test,alpha,beta))**2)
    sq_error_train_mean = np.mean((y_train - predict(X_train,alpha,beta))**2)
    error_test_4000.append(sq_error_test_mean)
    error_train_4000.append(sq_error_train_mean)

```

```

error_train_2000 = []
error_test_2000 = []
for i in lambda_list:
    alpha, beta = Neural_backprop(2000, 10, X_train, y_train, 0.01,i)
    sq_error_test_mean = np.mean((y_test - predict(X_test,alpha,beta))**2)
    sq_error_train_mean = np.mean((y_train - predict(X_train,alpha,beta))**2)
    error_test_2000.append(sq_error_test_mean)
    error_train_2000.append(sq_error_train_mean)

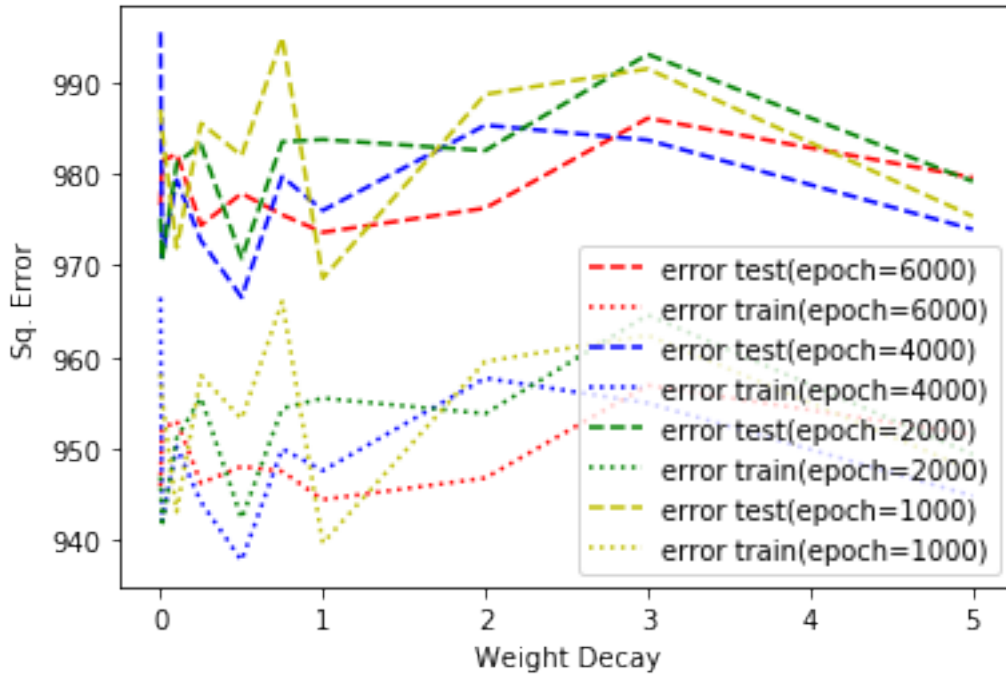
error_train_1000 = []
error_test_1000 = []
for i in lambda_list:
    alpha, beta = Neural_backprop(1000, 10, X_train, y_train, 0.01,i)
    sq_error_test_mean = np.mean((y_test - predict(X_test,alpha,beta))**2)
    sq_error_train_mean = np.mean((y_train - predict(X_train,alpha,beta))**2)
    error_test_1000.append(sq_error_test_mean)
    error_train_1000.append(sq_error_train_mean)

```

```

[58]: import matplotlib.pyplot as plt
fig, ax = plt.subplots()
plt.plot(lambda_list,error_test_6000,color='r', ls='--', label='error_
↳test(epoch=6000)')
plt.plot(lambda_list,error_train_6000,color='r', ls=':', label='error_
↳train(epoch=6000)')
plt.plot(lambda_list,error_test_4000,color='b', ls='--', label='error_
↳test(epoch=4000)')
plt.plot(lambda_list,error_train_4000,color='b', ls=':', label='error_
↳train(epoch=4000)')
plt.plot(lambda_list,error_test_2000,color='g',ls='--', label='error_
↳test(epoch=2000)')
plt.plot(lambda_list,error_train_2000,color='g',ls=':', label='error_
↳train(epoch=2000)')
plt.plot(lambda_list,error_test_1000,color='y',ls='--', label='error_
↳test(epoch=1000)')
plt.plot(lambda_list,error_train_1000,color='y',ls=':', label='error_
↳train(epoch=1000)')
plt.xlabel('Weight Decay')
plt.ylabel("Sq. Error")
leg = ax.legend();

```

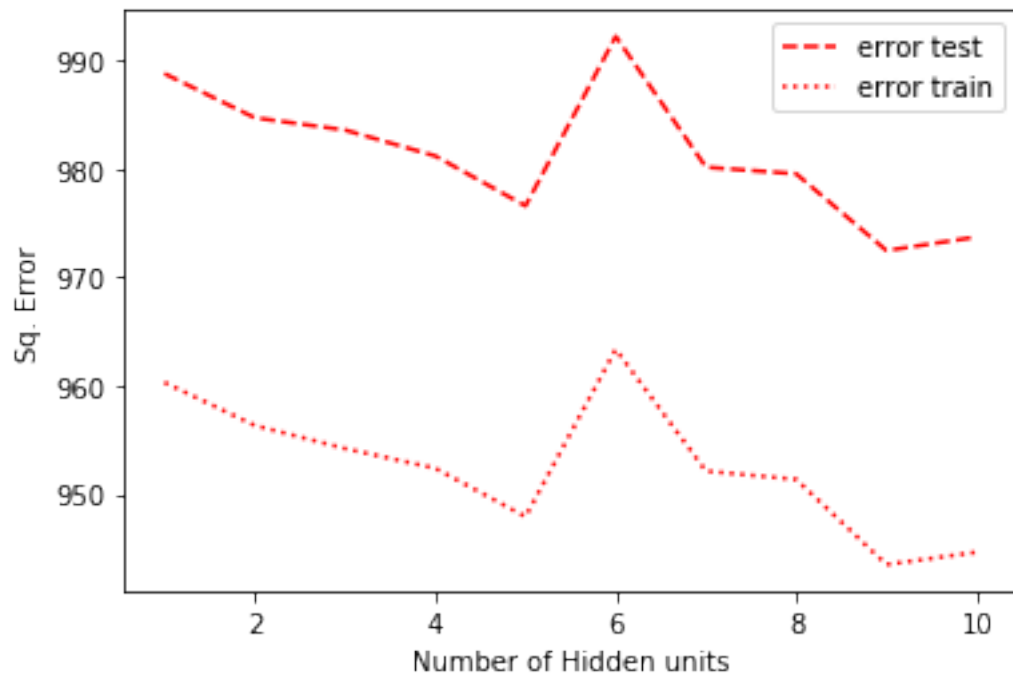


We can see that as increase the weight decay have different behaviours depending how the amount of epochs we pass. For example between 0 and 1 we can see that very sharp movements for the train and test error levels for the simulations with lower number epochs which means that are more susceptible to overfitting.

1.3 Part C

```
[61]: error_train_hn = []
error_test_hn = []
for i in hiddenSize_list:
    alpha, beta = Neural_backprop(6000, i, X_train, y_train, 0.01,.5)
    sq_error_test_mean = np.mean((y_test - predict(X_test,alpha,beta))**2)
    sq_error_train_mean = np.mean((y_train - predict(X_train,alpha,beta))**2)
    error_test_hn.append(sq_error_test_mean)
    error_train_hn.append(sq_error_train_mean)
```

```
[62]: import matplotlib.pyplot as plt
fig, ax = plt.subplots()
plt.plot(hiddenSize_list,error_test_hn,color='r', ls='--' ,label='error test')
plt.plot(hiddenSize_list,error_train_hn,color='r', ls=':', label='error train')
plt.xlabel('Number of Hidden units')
plt.ylabel("Sq. Error")
leg = ax.legend();
```



We can see that minimize at 9.