# HW3

1. Q1

Q1   A representation of the multilayered n/w

$x_1$ —— ○   $W_{11}^1$   $W_{01}^1$
     $W_{12}^1$
     $W_{13}$
$x_2$ —— ○
$x_3$ —— ○
     $W_{41}^1$
$x_4$ —— ○          Layer 1        Layer 2

This can be simplified to be represented
   as follows

layer (l-1)        layer (l)           layer (l+1)

○   $W_{ij}^{l-1}$        $z_i^l$         $x_i^l$        $z_m^{l+1}$
$x_j^{l-1}$        ○   (+)        ○   $W_{mi}^{l+1}$   ○

input at         weight at (l-1)          o/p at $l^{th}$
l-1 layer        layer between            layer
                 i,j units

- From here we get that input at layer l
  is the output of layer (l-1) & n/w feeds forward

Per the gradient descent update

$$W^{t+1} = W^t - \eta \, G(w^t)$$

From the diagram we can infer

$$z_i^\ell = w_{ij}^{(\ell-1)} \cdot x_j^{(\ell-1)} \quad —\ ①$$

hence $\dfrac{\partial z_i^\ell}{\partial w_{ij}} = (x_j^{l-1}) \quad —\ ②$

Also the assumption here is that

the transfer function $= \phi(z) = \dfrac{1}{1+e^{-z}} \quad —\ ③$

Let the output of the unit be $f_n(x)$

Error $= E_{in} = \dfrac{1}{2}[y_{in} - f_n(x_i)]^2$

For the $\ell^{th}$ layer

$$\frac{\partial E_{in}}{\partial z_i} = \frac{\partial}{\partial z_i} \frac{1}{2}[y_{in} - f_n(x_i)]^2$$

$$= \frac{1}{2} \frac{\partial}{\partial z_i}[y_{in} - \phi(z_i)]^2$$

$$= [y_{in} - \phi(z_i)]\left[\frac{\partial y_{in}}{\partial z_i} - \frac{\partial \phi(z_i)}{\partial z_i}\right]$$

if $\phi(z) = \dfrac{1}{1+e^{-z}}$ then $\phi'(z) = \dfrac{1}{(1+e^{-z})} \times \dfrac{e^{-z}}{(1+e^{-z})} = \phi(z)(1-\phi(z))$

hence

$$\frac{\partial E_{in}}{\partial z_i} = -[y_{in} - \phi(z_i)][\phi(z_i)][1-\phi(z_i)] \quad —\ ④$$

$$\frac{\partial E_{in}}{\partial z_i} = -\delta_{in} \qquad \longleftarrow \qquad ⑤$$

$$\text{when } \delta_{in} = -[y_{in} - \phi(z_i)][\phi(z_i)][1 - \phi(z_i)]$$

Now

$$w^{t+1} = w^t - \eta \, G(w^t) \qquad \longrightarrow \qquad ⑥$$

from ⑤ & ⑥ & ②

$$w^{t+1} = w^t \dotplus \eta \left[ \frac{\partial E_{in}}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}} \right]$$

$$w^{t+1} = w^t + \eta \, \delta_{in} \cdot x_{ij} \qquad \longrightarrow \qquad ⑦$$

Now consider the $(t+1)^{th}$ layer. The updates can be calculated using chain rule

$$\frac{\partial E_{in}}{\partial z_m} = \frac{\partial}{\partial z_m} [y_{in} - \phi(z_b)]^2$$

$$= -[y_{in} - \phi(z_i)] \cdot \frac{\partial \phi(z_m)}{\partial z_m}$$

$$= (1 - \phi(z_m))(\phi(z_m) \cdot \delta_{in} \, w_{ij} \qquad \text{from } ⑤$$

$$\frac{\delta E_{in}}{\delta z_m} \simeq \delta_{mn} \qquad \longrightarrow \qquad ⑧$$

$$\frac{\partial E_{in}}{\partial z_i} = \delta_{in}$$

from ⑤ & ⑧

$$\frac{\partial E_{in}}{\partial z_m} = \delta_{im}$$

We can continue to apply chain rule for each hidden layer to get δ's q the same form as above

Referring back to the diagram

① Calculate the feed forward $\phi(z_i)$

② Calculate the δ coefficient

$$\delta_{in} = \phi(z_i)(1-\phi(z_i)) \sum \delta_{jn} w_{ji}$$

(for all j's that are in layers downstream)

The update then is

$$w^{t+1} = w^t + \eta \, \delta_{in} \lambda_{ij}$$

which can be recursively calculated

2. Q2

(Q2)   $Q(w) = \sum_{i=1}^{N} q_i(w)$   where   $q_i(w) =$ error vector

$q_i(w) = \frac{1}{2}[y_i - F(x_i)]^2$

$F(x_i) = a_0 + \sum_{m=1}^{M} a_m B(\underline{x} | u_m, \sigma_m)$

$\quad = a_0 + \sum_{m=1}^{M} a_m \, e^{-\frac{1}{2\sigma_m^2} \sum_{i=1}^{n}(x_i - u_{im})^2}]$

Parameters are   $a_m$, $\sigma_m$ & $u_m$

gradient w.r.t $a_m$

$-\dfrac{\partial q_i(w)}{\partial a_m} = -\dfrac{\partial}{\partial a_m} \frac{1}{2}\left[y_i - a_0 - \sum_{m=1}^{M} a_m \, e^{-\frac{1}{2\sigma_m^2}\{(x_i - u_{im})^2}\right]$

$\quad = -[y_i - F(x_i)] \cdot \dfrac{\partial}{\partial a_m}\left[y_i - a_0 - \sum_{m=1}^{M} a_m \, e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - u_{im})^2}\right]$

$\quad = -[y_i - F(x_i)]\left[0 - 0 - \sum_{m=1}^{M} e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - u_{im})^2}\right]$

$\quad = -[y_i - F(x_i)] \, B(\underline{x} | u_m, \sigma_m)$  ⸺ ①

gradient w.r.t $U_m$

$$-\frac{\partial q_i(w)}{\partial U_m} = \frac{-1}{2} \cdot \frac{\partial}{\partial U_m}\left[ y_i - a_0 - \sum_{m=1}^{M} a_m e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_{im})^2} \right]$$

$$= -[y_i - F(x_i)] \cdot \frac{\partial}{\partial U_m}\left[ y_i - a_0 - \sum_{m=1}^{M} a_m e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_{im})^2} \right]$$

$$= -[y_i - F(x_i)]\left[ 0 - 0 - \frac{\partial}{\partial U_m}\sum_{m=1}^{M} a_m e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_{im})^2} \right]$$

$$= -[y_i - F(x_i)]\left[ a_m e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_m)^2} \right] \cdot \frac{\partial}{\partial U_m}\left[ -\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_{m})^2 \right]$$

$$= -[y_i - F(x_i)]\left[ a_m \cdot e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_{im})^2} \right]\left(-\frac{1}{2\sigma_m^2}\right)\left(2(x_i - U_{im})(-1)\right)$$

$$= -a_m[y_i - F(x_i)]\frac{\left[ e^{-\frac{1}{2\sigma_m^2}\sum_{i=1}^{n}(x_i - U_m)^2} \right][x_i - U_m]}{\sigma_m^2} \quad\quad ②$$

gradient w.r.t $\sigma_m$

$$-\frac{\partial q_i(w)}{\partial \sigma_m} = \frac{-1}{2} \frac{\partial}{\partial \sigma_m} \left[ y_i - a_0 - a_m \sum_{m=1}^{M} a_m e^{-\frac{1}{2\sigma_m^2} \sum_{l=1}^{n} (x_i - u_m)^2} \right]$$

$$= -\frac{2}{2} [y_i - F(x_i)] \left[ 0 - 6 - \frac{\partial}{\partial \sigma_m} \sum_{m=1}^{M} a_m e^{-\frac{1}{2\sigma_m^2} \sum_{l=1}^{n} (x_i - u_{lm})^2} \right]$$

$$= -[y_i - F(x_i)] \left[ -a_m e^{-\frac{1}{2\sigma_m^2} \sum_{l=1}^{n} (x_i - u_i)^2} \cdot \frac{\partial}{\partial \sigma_m} \left( -\frac{1}{2\sigma_m^2} \sum_{l=1}^{n} (x_i - u_{lm})^2 \right) \right]$$

$$= -a_m [y_i - F(x_i)] \left[ e^{-\frac{1}{2\sigma_m^2} \sum_{l=1}^{n} (x_i - u_{lm})^2} \cdot \sum (x_i - u_{lm})^2 \frac{(-1)(-2)}{2 \sigma_m^3} \right]$$

$$= -a_m [y_i - F(x_i)] \left[ \frac{e^{-\frac{1}{2\sigma_m^2} \sum_{l=1}^{n} (x_i - u_{lm})^2}}{\sigma_m^3} \right] \left[ \sum_{l=1}^{n} (x_i - u_m)^2 \right] \quad -③$$

①, ②, ③ show the derivative of gradient

## 3. Q3

Q3

$$B(X \mid U_m, \Sigma) = e^{-\frac{1}{2}\left[(x-u_m)^T \Sigma (x-u_m)\right]}$$

$\Sigma =$ positive definite and symmetric matrix

hence $\quad \Sigma = \Sigma^{1/2} \Sigma^{1/2} \quad\quad - \text{(1)}$

Consider an input vector $x$ that is pre-multiplied
with $\Sigma^{1/2}$

i.e $\quad u/p = \Sigma^{-1/2} x = \hat{x}$

$$B(\hat{x} \mid \hat{u}_m, \Sigma) = e^{-\frac{1}{2}\left[\Sigma^{1/2}(\hat{x}-\hat{u}_m)^T \Sigma \Sigma^{1/2}(\hat{x}-\hat{u}_m)\right]}$$

$$= e^{-\frac{1}{2}\left[(\hat{x}-\hat{u}_m) \cdot \Sigma^{1/2} \cdot \Sigma \cdot \Sigma^{-1/2}(\hat{x}-u_m)\right]}$$

using (1)

$$\doteq e^{-\frac{1}{2}\left[(\hat{x}-\hat{u}_m)^T [I] (\hat{x}-\hat{u}_m)\right]}$$

$$B(\hat{x} \mid \hat{u}_m, \Sigma) = e^{-\frac{1}{2}\left[(\hat{x}-\hat{u}_m)^T (x-u_m)\right]} \quad - \text{(2)}$$

In problem (2) if $\sigma_m = 1$

then $\quad B(x \mid u_m, \sigma_m) = e^{-\frac{1}{2}\left[x-u_m\right]^2} \quad - \text{(3)}$

then (2) & (3) are equivalent
hence proved

4. Q4

Q4

$$B(\underline{x} \mid u_m, \Sigma_m) = e^{-\frac{1}{2}\left[(x-u_m)^t \Sigma_m (x-u_m)\right]}$$

→ In the std neural n/w the transfer function vary only in 1 direction of i/p space

→ Radial bases functions, vary in all directions equally

→ $\Sigma_m$ controls the direction of variation

→ If there exists a vector $\bar{x}$ such that
$$\Sigma_m \bar{x} = \lambda \bar{x}$$

ie if the $\bar{x} =$ eigen vector of $\Sigma_m$ then the direction of variedness is just along that vector direction

5. Q5.

    – 5a.

```r
library(nnet)

inspam=read.csv("Spam_Train.txt")
spname<-c ("make", "address", "all", "3d", "our", "over", "remove",
            "internet","order", "mail", "receive", "will",
            "people", "report", "addresses","free", "business",
            "email", "you", "credit", "your", "font","000","money",
            "hp", "hpl", "george", "650", "lab", "labs",
            "telnet", "857", "data", "415", "85", "technology", "1999",
            "parts","pm", "direct", "cs", "meeting", "original", "project",
            "re","edu", "table", "conference", ";", "(", "[", "!", "$", "#",
            "CAPAVE", "CAPMAX", "CAPTOT","type")
colnames(inspam)=spname
x=inspam
colnames(x)=spname
x$type=as.factor(inspam$type)
x[,1:57]=scale(x[,1:57], center=TRUE, scale=TRUE)
x=data.frame(x)
colnames(x)=spname


inspamtest=read.csv("Spam.Test.txt")
colnames(inspamtest)=spname
w=inspamtest
colnames(w)=spname
w$type=as.factor(as.factor(inspamtest$type))
w[,1:57]=scale(w[,1:57],center=TRUE, scale=TRUE)
w=data.frame(w)
colnames(w)=spname
w.type=w$type

#Find the size with minimum error
set.seed(1)
for(i in 1:10){
  nn1=nnet(type~.,data=x, size=i, maxit=5000, decay=0.0, rang=0.5, trace=F)
  nn1.predict=predict(nn1, newdata = w,type="class")
  nn1.out=nn1.predict
  u=matrix(data=0,2,2)
  u=table(nn1.out, w.type)
  err=sum(nn1.out!=w.type)/(length(nn1.out))
  print(paste0("# of hidden nodes = ",i," and error = ", err))
}

## [1] "# of hidden nodes = 1 and error = 0.0730593607305936"
## [1] "# of hidden nodes = 2 and error = 0.0743639921722113"
## [1] "# of hidden nodes = 3 and error = 0.0652315720808871"
## [1] "# of hidden nodes = 4 and error = 0.060665362035225"
## [1] "# of hidden nodes = 5 and error = 0.0639269406392694"
## [1] "# of hidden nodes = 6 and error = 0.0574037834311807"
```

```
## [1] "# of hidden nodes = 7 and error = 0.0600130463144162"
## [1] "# of hidden nodes = 8 and error = 0.0547945205479452"
## [1] "# of hidden nodes = 9 and error = 0.0678408349641226"
## [1] "# of hidden nodes = 10 and error = 0.108936725375082"

print(paste0("Minimum Error is with # of hidden nodes = ", 8))

## [1] "Minimum Error is with # of hidden nodes = 8"
```

+ By using a # of hidden nodes as 8, we get  overall error rate ~4-5%

- 5b.

```
set.seed(1)
res=matrix(NA, length(nn1.out),11)
ii=1;
for(j in seq(0,1,0.1)){
  nn1=nnet(type~.,data=x, size=8, maxit=5000, decay=j, rang=-0.5, trace=F)
  nn1.predict=predict(nn1, newdata = w[,1:57],type="class")
  nn1.out=nn1.predict
  res[,ii]=nn1.out
  ii=ii+1;
  err=sum(nn1.out!=w.type)/(length(nn1.out))
  print(paste0("Decay = ",j," and error = ", err))
}

## [1] "Decay = 0 and error = 0.065883887801696"
## [1] "Decay = 0.1 and error = 0.0482713633398565"
## [1] "Decay = 0.2 and error = 0.0547945205479452"
## [1] "Decay = 0.3 and error = 0.0430528375733855"
## [1] "Decay = 0.4 and error = 0.0437051532941944"
## [1] "Decay = 0.5 and error = 0.0476190476190476"
## [1] "Decay = 0.6 and error = 0.0521852576647097"
## [1] "Decay = 0.7 and error = 0.0476190476190476"
## [1] "Decay = 0.8 and error = 0.0528375733855186"
## [1] "Decay = 0.9 and error = 0.0476190476190476"
## [1] "Decay = 1 and error = 0.0547945205479452"

print(paste0("Minimum Error is with value of decay as = ", 0.6))

## [1] "Minimum Error is with value of decay as = 0.6"

nn1=nnet(type~.,data=x, size=8, maxit=5000, decay=0.6, rang=-0.5, trace=F)
nn1.best=predict(nn1, newdata = w[,1:57],type="class")

#Finding the class through majority vote:
vote=rep(NA,length(nn1.out))
for (i in 1:length(nn1.out)){
  if(sum(res[i,]==1)>sum(res[i,]==0))
  {vote[i]=1}
  else{
    vote[i]=0
  }
}

#Calcualte error
err=sum(vote!=w.type)/(length(nn1.out))
print(paste0("Error using the majority of votes is  ", err))

## [1] "Error using the majority of votes is  0.0430528375733855"
```

```
  + Best Model:
    + Decay: 0.6
    + Number of hidden units: 8
    + Error:~3-4%


  + By using an ensemble, where we find majority of votes for a class, the error is about
4%
```

- 5c.

```r
set.seed(1)
for(k in seq(0,1,0.1)){
    nn1=nnet(type~.,data=x, size=8, maxit=5000, decay=0.6, trace=F)
    nn1.predict=predict(nn1, newdata = w,type="raw")
    nn1.out=rep(0,length(nn1.predict))
    nn1.out[nn1.predict>k]=1
    u=matrix(data=0,2,2)
    u=table(w.type,nn1.out)

    print(paste0("Threshold= ", k , " Proportion of good mails misclassified is: ", u[3
][1]/(u[1][1]+u[3][1]) ))
}

## [1] "Threshold= 0 Proportion of good mails misclassified is: 0.998908296943231"
## [1] "Threshold= 0.1 Proportion of good mails misclassified is: 0.140829694323144"
## [1] "Threshold= 0.2 Proportion of good mails misclassified is: 0.0927947598253275"
## [1] "Threshold= 0.3 Proportion of good mails misclassified is: 0.0676855895196507"
## [1] "Threshold= 0.4 Proportion of good mails misclassified is: 0.0524017467248908"
## [1] "Threshold= 0.5 Proportion of good mails misclassified is: 0.0305676855895196"
## [1] "Threshold= 0.6 Proportion of good mails misclassified is: 0.0316593886462882"
## [1] "Threshold= 0.7 Proportion of good mails misclassified is: 0.0196506550218341"
## [1] "Threshold= 0.8 Proportion of good mails misclassified is: 0.0163755458515284"
## [1] "Threshold= 0.9 Proportion of good mails misclassified is: 0.00655021834061135"
## [1] "Threshold= 1 Proportion of good mails misclassified is: NA"

  + By using a Threshold of 90% of classifiying an email as spam, we can get misclassific
ation of good email as spam down to <1% error rate.
```