

1.

$$1. \quad \theta(s) = \frac{e^s}{1+e^s} = \frac{1}{1+e^{-s}}$$

$$\text{假设 } F(A, B) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n(Az_n + B)})$$

$$\boxed{p_n = \theta(-y_n(Az_n + B))} \\ = \frac{1}{1 + e^{y_n(Az_n + B)}}$$

$$= -\frac{1}{N} \sum_{n=1}^N \ln\left(\frac{1}{1+e^{-s}}\right) = -\frac{1}{N} \sum_{n=1}^N \ln\left(\frac{e^s}{1+e^s}\right) = -\frac{1}{N} \sum_{n=1}^N \ln(1-p_n)$$

$$\frac{\partial p_n}{\partial A} = -p_n(1-p_n)(y_n z_n)$$

$$\frac{\partial p_n}{\partial B} = -p_n(1-p_n)(y_n)$$

$$\frac{\partial F}{\partial A} = -\frac{1}{N} \sum_{n=1}^N \frac{1}{1-p_n} (-1) \frac{\partial p_n}{\partial A}$$

$$\frac{\partial F}{\partial B} = -\frac{1}{N} \sum_{n=1}^N \frac{1}{1-p_n} (-1) \frac{\partial p_n}{\partial B}$$

$$= -\frac{1}{N} \sum_{n=1}^N \frac{1}{1-p_n} (-1) (-p_n)(1-p_n)(y_n z_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N \frac{1}{1-p_n} (-1) (-p_n)(1-p_n)(y_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N p_n y_n z_n$$

$$= -\frac{1}{N} \sum_{n=1}^N p_n y_n$$

$$\nabla F(A, B) = \left[ -\frac{1}{N} \sum_{n=1}^N p_n y_n z_n, -\frac{1}{N} \sum_{n=1}^N p_n y_n \right]$$

2.

2.

$$H(F) = \begin{bmatrix} \frac{\partial^2 F(A,B)}{\partial A^2} & \frac{\partial^2 F(A,B)}{\partial A \partial B} \\ \frac{\partial^2 F(A,B)}{\partial A \partial B} & \frac{\partial^2 F(A,B)}{\partial B^2} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N z_n^2 p_n (1-p_n) & \frac{1}{N} \sum_{n=1}^N z_n p_n (1-p_n) \\ \frac{1}{N} \sum_{n=1}^N z_n p_n (1-p_n) & \frac{1}{N} \sum_{n=1}^N p_n (1-p_n) \end{bmatrix}$$

$$\frac{\partial^2 F(A,B)}{\partial A^2} = \frac{\partial}{\partial A} \left( -\frac{1}{N} \sum_{n=1}^N p_n y_n z_n \right)$$

$$= -\frac{1}{N} \sum_{n=1}^N y_n z_n \frac{\partial p_n}{\partial A} = -\frac{1}{N} \sum_{n=1}^N y_n z_n (-p_n)(1-p_n)(y_n z_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N (y_n z_n)^2 (-p_n)(1-p_n) = \frac{1}{N} \sum_{n=1}^N z_n^2 p_n (1-p_n)$$

$$\frac{\partial^2 F(A,B)}{\partial B^2} = \frac{\partial}{\partial B} \left( -\frac{1}{N} \sum_{n=1}^N p_n y_n \right)$$

$$= -\frac{1}{N} \sum_{n=1}^N y_n \frac{\partial p_n}{\partial B} = -\frac{1}{N} \sum_{n=1}^N y_n (-p_n)(1-p_n)y_n = \frac{1}{N} \sum_{n=1}^N p_n (1-p_n)$$

$$\frac{\partial^2 F(A,B)}{\partial A \partial B} = \frac{\partial}{\partial B} \left( -\frac{1}{N} \sum_{n=1}^N p_n y_n z_n \right)$$

$$= -\frac{1}{N} \sum_{n=1}^N y_n z_n \frac{\partial p_n}{\partial B} = -\frac{1}{N} \sum_{n=1}^N y_n z_n (-p_n)(1-p_n)y_n = \frac{1}{N} \sum_{n=1}^N z_n (p_n)(1-p_n)$$

3.

3. We can prove that the matrix  $H(F)$  is positive semi-definite by finding all of its eigenvalue  $\geq 0$ .

Let the eigenvalue  $= \lambda$

$$H(F)B = \lambda B = (\lambda I)B$$

$$(H(F) - \lambda I)B = 0$$

$$H(F) - \lambda I = \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N z_n^2 P_n (1 - P_n) - \lambda & \frac{1}{N} \sum_{n=1}^N z_n P_n (1 - P_n) \\ \frac{1}{N} \sum_{n=1}^N z_n P_n (1 - P_n) & \frac{1}{N} \sum_{n=1}^N P_n (1 - P_n) - \lambda \end{bmatrix}$$

Let  $\frac{1}{N} \sum_{n=1}^N P_n (1 - P_n) = \alpha$

$$\det(H(F) - \lambda I) = 0$$

$$= (z_n^2 \alpha - \lambda)(\alpha - \lambda) - (z_n \alpha)(z_n \alpha)$$

$$= \cancel{z_n^2 \alpha^2} - \cancel{z_n^2 \alpha^2} + \lambda^2 - \lambda(\alpha + z_n^2 \alpha)$$

$$= \lambda(\lambda - \alpha - z_n^2 \alpha) = 0$$

$\lambda = 0 (\geq 0) \Rightarrow \text{ok}$  between 0 and 1

$\lambda = \alpha + z_n^2 \alpha = \alpha(1 + z_n^2) = \frac{1}{N} \sum_{n=1}^N P_n (1 - P_n) (1 + z_n^2) (\geq 0) \Rightarrow \text{ok}$   
 between 0 and 1 positive

With the eigenvalues 0 and  $\alpha + z_n^2 \alpha \geq 0$ , we can prove that  $H(F)$  is positive semi-definite.

4.

4. 1st sign(0) = +1

$(w_0, w_1, w_2, \dots, w_d) = (d-1, +1, +1, \dots, +1)$

$$gA(x) = \text{sign} \left( \sum_{i=1}^d x_i + d - 1 \right)$$

Only when all the  $x_n$  is -1,  $gA(x) = \text{sign}(-d + 1 - 1) = -1$ , which means that if any  $x_n = +1$ ,  $gA(x)$  will be 1.

5.

5.

$$\frac{\partial e_n}{\partial w_{ij}^{(l)}} = \delta_j^{(l)} (\gamma_i^{(l-1)})$$

$$\delta_j^{(l)} = \frac{\partial e_n}{\partial z_j^{(l)}} = \sum_k (\delta_k^{(l+1)}) (w_{jk}^{(l+1)}) (\tanh'(z_j^{(l)}))$$

$$\gamma_i^{(l)} = 0 \text{ for } l \geq 1$$

$$\delta_j^{(l)} = 0 \text{ for } l < 6, \delta_j^{(6)} \neq 0$$

if  $y_n \neq 0$ ,  $\gamma_i^{(6-1)} = 0$ ,  $z \neq 0$  and  $\gamma_0^{(6-1)} = 1$

All the gradient components must be zero except for the gradient components with respect to  $w_{01}^{(1)}$

6.

b.

$$\underbrace{d^{(0)}}_{\text{input}} + \underbrace{1}_{\text{constant}} = 12 \quad \sum_{n=1}^N \gamma_n = 48$$

max:  $12 \cdot \gamma_1 + \gamma_1 \cdot \gamma_2 + \dots + \gamma_n \cdot 1$

while  $N=1$ ,  $12 \times (48-1) + 48 \times 1 = 612$

while  $N=2$ ,  $12 \times (\gamma_1 - 1) + (\gamma_1)(48 - \gamma_1 - 1) + (48 - \gamma_1)$

$$= 12\gamma_1 - 12 + 47\gamma_1 - \gamma_1^2 + 48 - \gamma_1$$

$$= -\gamma_1^2 + 58\gamma_1 + 36 = -(\gamma_1 - 29)^2 + 877$$

NNet = 11 - 28 - 18 - 1

(11+1) (28+1) (18+1)

max:  $12 \times 28 + 29 \times 18 + 19 \times 1 = 336 + 522 + 19 = 877$



7.

$$\begin{aligned}
 \text{err}_n(w) &= \|\chi_n - ww^T \chi_n\|^2 \\
 &= (\chi_n - ww^T \chi_n)^T (\chi_n - ww^T \chi_n) \\
 &= \chi_n^T \chi_n - 2(\chi_n^T w)(w^T \chi_n) + \chi_n^T ww^T ww^T \chi_n \\
 &= \chi_n^T \chi_n - 2 \underbrace{(w^T \chi_n)^2}_{(-2w^T \chi_n)(w^T \chi_n)} + (w^T \chi_n)^2 (w^T w) \\
 \nabla_w \text{err}(w) &= -2\chi_n(w^T \chi_n) - 2w^T \chi_n(\chi_n) + \frac{\partial}{\partial w} (w^T \chi_n)^2 (w^T w) \\
 &= -4w^T \chi_n \chi_n + (w^T \chi_n)^2 2w + (w^T w) \cdot \frac{\partial}{\partial w} (w^T \chi_n)^2 \\
 &= -4w^T \chi_n \chi_n + 2(w^T \chi_n)^2 w + (w^T w)(w^T \chi_n \chi_n + w^T \chi_n \chi_n) \\
 &= \underbrace{-4w^T \chi_n \chi_n}_{d \times 1} + \underbrace{2(w^T \chi_n)^2 w}_{d \times 1} + \underbrace{2(w^T w)(w^T \chi_n \chi_n)}_{d \times 1}
 \end{aligned}$$

8.

$$\begin{aligned}
 E_{in}(w) &= \frac{1}{N} \sum_{n=1}^N \|\chi_n - ww^T (\chi_n + \epsilon_n)\|^2 \\
 &= \frac{1}{N} \sum_{n=1}^N (\chi_n^T \chi_n - 2\chi_n^T ww^T (\chi_n + \epsilon_n) + (\chi_n + \epsilon_n)^T ww^T ww^T (\chi_n + \epsilon_n)) \\
 &= \frac{1}{N} \sum_{n=1}^N (\chi_n^T \chi_n - 2\chi_n^T ww^T \chi_n - 2\chi_n^T ww^T \epsilon_n + \chi_n^T ww^T ww^T \chi_n + \epsilon_n^T ww^T ww^T \epsilon_n + 2\epsilon_n^T ww^T \chi_n) \\
 &= \frac{1}{N} \sum_{n=1}^N (\chi_n^T \chi_n - 2\chi_n^T ww^T \chi_n + \chi_n^T ww^T ww^T \chi_n) + \frac{1}{N} \sum_{n=1}^N (-2\chi_n^T ww^T \epsilon_n + \epsilon_n^T ww^T ww^T \epsilon_n + 2\epsilon_n^T ww^T \chi_n) \\
 &= \frac{1}{N} \sum_{n=1}^N \|\chi_n - ww^T \chi_n\|^2 + \frac{1}{N} \sum_{n=1}^N (-2\chi_n^T ww^T \epsilon_n + \epsilon_n^T ww^T ww^T \epsilon_n + 2\epsilon_n^T ww^T \chi_n)
 \end{aligned}$$

We know that  $E[\epsilon_n \epsilon_n^T] = I_n$   
 and the  $\epsilon_n$ 's first order's Expectation = 0

$$\begin{aligned}
 E[E_{in}(w)] &= \frac{1}{N} \sum_{n=1}^N \|\chi_n - ww^T \chi_n\|^2 + \frac{1}{N} \sum_{n=1}^N E[\epsilon_n^T ww^T ww^T \epsilon_n] \\
 &= \frac{1}{N} \sum_{n=1}^N \|\chi_n - ww^T \chi_n\|^2 + \frac{1}{N} \sum_{n=1}^N E[w^T \epsilon_n \epsilon_n^T w] \\
 &= \frac{1}{N} \sum_{n=1}^N \|\chi_n - ww^T \chi_n\|^2 + \frac{1}{N} \sum_{n=1}^N w^T w E[\epsilon_n^T \epsilon_n] w \\
 &= \frac{1}{N} \sum_{n=1}^N \|\chi_n - ww^T \chi_n\|^2 + (w^T w)^2 \\
 &= \Omega(w) = (w^T w)^2
 \end{aligned}$$

9.

error function  $\sum_{i=1}^d (g_i(x) - x_i)^2$

$$= \sum_{i=1}^d (x_i w_{ij}^{(1)} w_{ji}^{(2)} - x_i)^2$$

$$= \sum_{i=1}^d (x_i u_{ij} u_{ij}^T - x_i)^2$$

10.

10.

$$\frac{\partial E_d(u)}{\partial u_{ij}} = \frac{\partial E_o(w)}{\partial w_{ij}^{(1)}} + \frac{\partial E_o(w)}{\partial w_{ij}^{(2)}}$$

$$= \sum_{i=1}^d \frac{\partial}{\partial u_{ij}} (x_i u_{ij} u_{ij}^T - x_i)^2 = \sum_{i=1}^d 2(x_i u_{ij} u_{ij}^T - x_i)(x_i u_{ij})$$

$$= \sum_{i=1}^d 2(x_i u_{ij} u_{ij}^T - x_i)(x_i u_{ij}) = \sum_{i=1}^d 2(x_i w_{ij}^{(1)} w_{ji}^{(2)} - x_i)(x_i w_{ij}^{(1)})$$

$$= \sum_{i=1}^d 2(x_i w_{ij}^{(1)} w_{ji}^{(2)} - x_i)(x_i w_{ij}^{(1)}) = \sum_{i=1}^d 2(x_i w_{ij}^{(1)} w_{ji}^{(2)} - x_i)(x_i w_{ij}^{(1)})$$

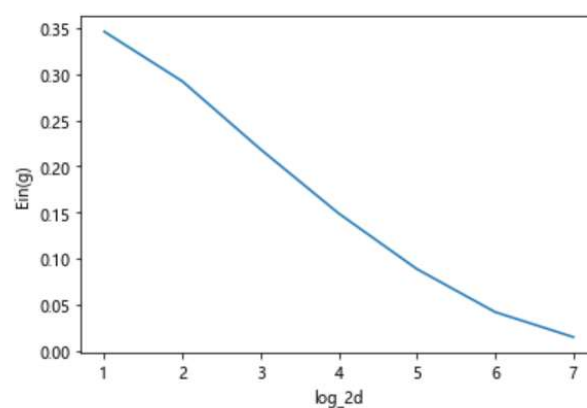
We know that  $u_{ij} = w_{ij}^{(1)}$ , and  $u_{ij} = w_{ji}^{(2)}$

So we can prove that

$$\frac{\partial E_d(u)}{\partial u_{ij}} = \frac{\partial E_o(w)}{\partial w_{ij}^{(1)}} + \frac{\partial E_o(w)}{\partial w_{ji}^{(2)}}$$

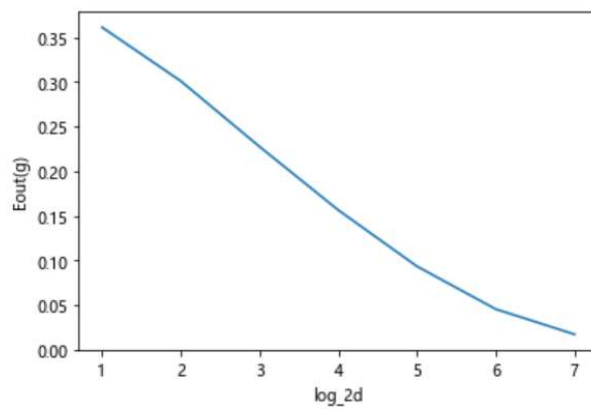
11.

Out[93]: Text(0, 0.5, 'Ein(g)')



As d goes bigger, the size of hidden layer gets bigger, the Ein(g) gets smaller.

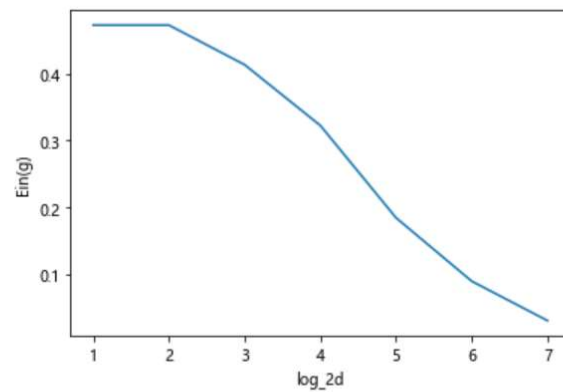
12.



Same with 11, as  $d$  goes bigger, the size of hidden layer gets bigger, the  $E_{in}(g)$  gets smaller. However,  $E_{out}$ 's error rate is bigger than  $E_{in}$  in every  $d$ , because the model is fit by training data but not testing data.

13.

Out[16]: Text(0, 0.5, 'E<sub>in</sub>(g)')

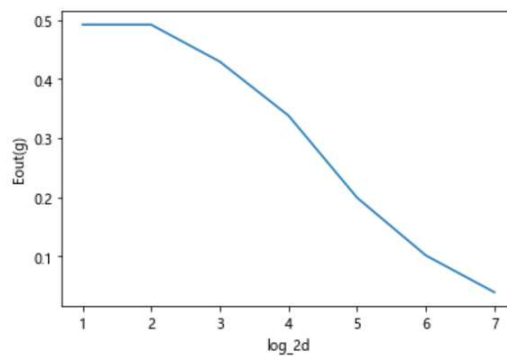


13.

As  $d$  goes bigger, the size of hidden layer goes bigger,  $E_{in}(g)$  gets smaller, but still bigger than autoencoder in 11.

14.

Out[17]: Text(0, 0.5, 'E<sub>out</sub>(g)')



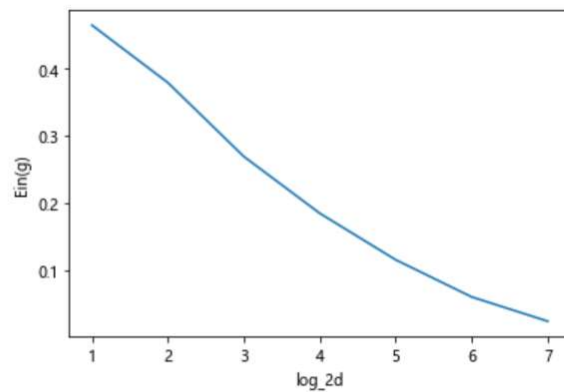
14.

Same with 13, as  $d$  goes bigger, the size of hidden layer goes bigger,  $E_{in}(g)$  gets smaller, but still bigger than autoencoder in 12. Also, the error rate that the test data performs are worse than 13.

15.

```
[0.46454766, 0.3790338, 0.26911533, 0.18459733, 0.11539804, 0.060169443, 0.02405729]
```

```
Out[12]: Text(0, 0.5, 'Ein(g)')
```



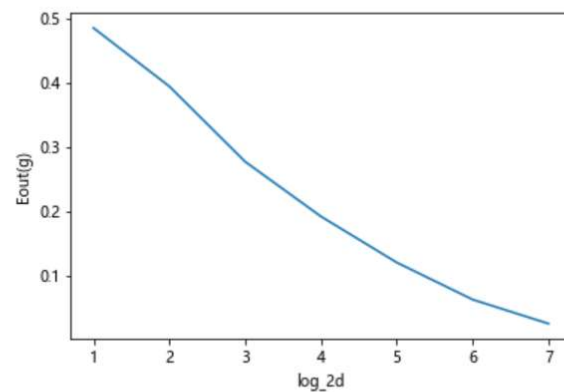
15.

As  $d$  goes bigger, the size of hidden layer goes bigger,  $\text{Ein}(g)$  gets smaller.  
Every  $\text{Ein}(g)$  in PCA perform better than  $\text{Ein}(g)$  in 13. in each different  $d$ .

16.

```
[0.4847186, 0.3939795, 0.2769158, 0.19177727, 0.12019751, 0.06265429, 0.025237955]
```

```
Out[14]: Text(0, 0.5, 'Eout(g)')
```



15.

As  $d$  goes bigger, the size of hidden layer goes bigger,  $\text{Ein}(g)$  gets smaller.  
Compares to  $\text{Ein}$ ,  $\text{Eout}$  is bigger(error is larger).  
Every  $\text{Ein}(g)$  in PCA perform better than  $\text{Ein}(g)$  in 14. in each different  $d$ .



17.

17.

We want to prove  $N^0 < 2^N$

First, we change it to  $\Delta \log N < N \log 2$  by using log.

So, we know  $f(N) = \Delta \log N - N \log 2 < 0$ .

$f'(N) = \frac{\Delta}{N} - \log 2$

while  $\frac{\Delta}{N} - \log 2 > 0$ ,  $f'(N) > 0$

$\Rightarrow \frac{\Delta}{N} > \log 2$

$\Rightarrow N < \frac{\Delta}{\log 2}$ ,  $f'(N) > 0$

while  $\frac{\Delta}{N} - \log 2 < 0$ ,  $f'(N) < 0$

$\Rightarrow N > \frac{\Delta}{\log 2}$ ,  $f'(N) < 0$

$N = 3 \Delta \log_2 \Delta = \frac{3 \Delta \ln \Delta}{\ln 2} > \frac{\Delta}{\ln 2}$  while  $\Delta \geq 2$ ,

so we can know that

$f(N) \leq f(3 \Delta \log_2 \Delta) < 0$

$f'(N) < 0$

Finally, we can know that while  $N \geq 3 \Delta \log_2 \Delta$ ,

$N^0 < 2^N$

and while  $\Delta, N$  are integers, we can prove that

$N^0 + 1 < 2^N$

18.

18.

First, we know that

$$B(N, k) \leq \sum_{i=0}^{k-1} C_i^N \quad \text{and} \quad \sum_{i=0}^D C_i^N \leq N^D + 1$$

From the question, we know that the hypothesis  $H_N$  that consists of all  $d-3-1$  neural networks' maximum combination amount is

$$B(N, d+1)^3$$

Because

$$B(N, d+1) \leq \sum_{i=0}^d C_i^N \leq N^d + 1 \leq N^{d+1} + 1, \quad \text{we can know that}$$

$$B(N, d+1)^3 = (N^{d+1} + 1)^3 = N^{3(d+1)} + 3N^{2(d+1)} + 3N^{d+1} + 1$$

$$\text{and while } \Delta = 3(d+1) + 1, \quad N \geq 0 \geq 4,$$

$$B(N, d+1)^3 = N^{3(d+1)} + 3N^{2(d+1)} + 3N^{d+1} + 1$$

$$< N^{3(d+1)} + N^{3(d+1)} + N^{3(d+1)} + 1 = 3N^{3(d+1)} + 1 < N^{3(d+1)+1} + 1$$

$$\text{while } N \geq 3 \geq \log_2 6,$$

$$(N^{d+1} + 1)^3 < N^{3(d+1)+1} + 1 \leq 2^N$$

$$VC < 3 \geq \log_2 \Delta = 3(3(d+1)+1) \log_2 (3(d+1)+1)$$