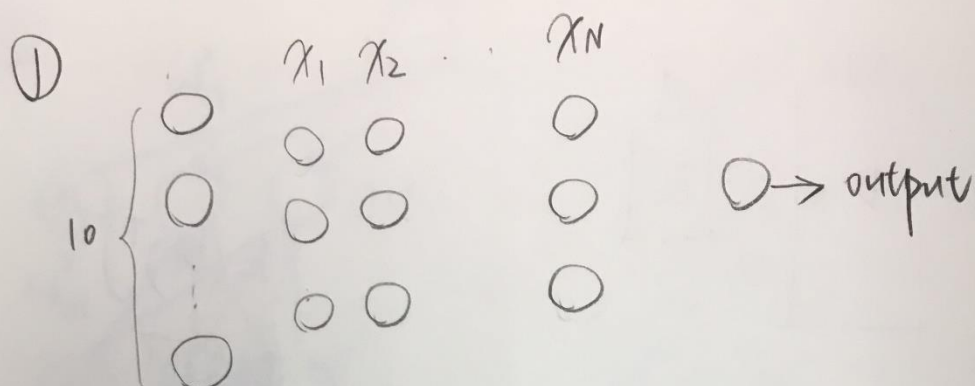


Problem 1



Let the i th layer has x_i neurons.

$$\sum_{n=1}^N x_n = 36, \quad x_n \geq 2 \quad \forall n=1, \dots, N$$

Total weights number

$$\begin{aligned} &= 10(x_1 - 1) + x_1(x_2 - 1) + \dots + x_N \cdot 1 \\ &= 10x_1 + x_1x_2 + \dots + x_N \cdot 1 - (10 + x_1 + \dots + x_{N-1}) \\ &= 10x_1 + x_1x_2 + \dots + x_N \cdot 1 - (10 + 36 - x_N) \\ &= (10x_1 + x_1x_2 + \dots + x_N \cdot 2) - 46 \end{aligned}$$

For each neuron to any layer \rightarrow at least has 4 weights
(product of neuron number of preceding & back layers).

$\Rightarrow x_1 = x_2 = \dots = x_N = 2$ leads to minimum weights

$$\Rightarrow (10 \cdot 2 + 4 \cdot 18) - 46 = 46 \text{ (minimum)}$$

#

Problem 2

②

According to question ①.

$$\text{Total weights number} = (10x_1 + x_1x_2 + \dots + x_N \cdot 2) - 46$$

$$\text{and } \sum_{i=1}^N x_i = 36, \quad x_i \geq 2 \quad \forall i = 1, 2, \dots, N$$

$$\text{Assume 1 hidden layer: } 10 \cdot 36 + 36 \cdot 2 - 46 = 386.$$

Assume 2 hidden layer:

$$\text{let } x_1 = x, \quad x_2 = 36 - x$$

$$10x + x(36 - x) + (36 - x) \cdot 2 - 46 = -x^2 + 44x + 26 = -(x - 22)^2 + 510$$

when $x = 22$ has max 510

Assume ≥ 3 hidden layer. $\Rightarrow < 510$ weights.

\Rightarrow maximum possible number of weights = 510 #

Problem 3

③

$$\text{err}_n(w) = \|x_n - ww^T x_n\|^2$$

$$= (x_n - ww^T x_n)^T (x_n - ww^T x_n)$$

$$= (x_n^T - x_n^T w w^T) (x_n - ww^T x_n)$$

$$= x_n^T x_n - 2x_n^T w w^T x_n + x_n^T w w^T w w^T x_n$$

$$= x_n^T x_n - 2(w^T x_n)^2 + (w^T x_n)^2 (w^T w)$$

$$(\because w^T x_n = x_n^T w = k, k \text{ is a constant})$$

$$\nabla_w \text{err}_n(w) = \frac{\partial x_n^T x_n}{\partial w} - 4(w^T x_n) \frac{\partial w^T x_n}{\partial w}$$

$$+ 2(w^T x_n) \frac{\partial w^T x_n}{\partial w} (w^T w) + (w^T x_n)^2 \frac{\partial w^T w}{\partial w}$$

$$= -4(w^T x_n)x_n + 2(w^T x_n)(w^T w)x_n + 2(w^T x_n)^2 w$$

#

Problem 4

④ $\epsilon_n \sim \text{uniform}(0, 1)$.

$$\begin{aligned}
 E_n(w) &= \frac{1}{N} \sum_{n=1}^N \|x_n - ww^T(x_n + \epsilon_n)\|^2 \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - ww^T x_n - ww^T \epsilon_n\|^2 \quad (\text{let } k_n = x_n - ww^T x_n) \\
 &= \frac{1}{N} \sum_{n=1}^N (k_n - ww^T \epsilon_n)^T (k_n - ww^T \epsilon_n) \\
 &= \frac{1}{N} \sum_{n=1}^N (k_n^T - \epsilon_n^T ww^T) (k_n - ww^T \epsilon_n) \\
 &= \frac{1}{N} \sum_{n=1}^N (k_n^T k_n - k_n^T ww^T \epsilon_n - \epsilon_n^T ww^T k_n + \epsilon_n^T ww^T ww^T \epsilon_n) \\
 &= \frac{1}{N} \sum_{n=1}^N [k_n^T k_n - 2(w^T k_n)(w^T \epsilon_n) + (w^T \epsilon_n)^2 (w^T w)] \\
 &= \frac{1}{N} \sum_{n=1}^N k_n^T k_n - 2 \underbrace{E[(w^T k_n)(w^T \epsilon_n)]}_0 + E[(w^T \epsilon_n)^2] \cdot (w^T w) \\
 &= \frac{1}{N} \sum_{n=1}^N k_n^T k_n - 0 + E[w^T \epsilon_n \epsilon_n^T w] \cdot (w^T w)
 \end{aligned}$$

and $\epsilon_{n,i} \sim \text{uniform}(0, 1)$.

$$E[\epsilon_{n,i} \cdot \epsilon_{n,k}] = \begin{cases} 0 & \text{if } i \neq k \quad (\because E[\epsilon_{n,i}] \cdot E[\epsilon_{n,k}] = 0) \\ 1 & \text{if } i = k \quad (\because E[x^2] = (E[x])^2 + \text{Var}(x) = 1) \end{cases}$$

$$\begin{aligned}
 \therefore E_n(w) &= \frac{1}{N} \sum_{n=1}^N k_n^T k_n + E[w^T I w] \cdot (w^T w) \\
 &= \frac{1}{N} \sum_{n=1}^N \|x_n - ww^T x_n\|^2 + (w^T w)^2 \\
 \therefore \Omega(w) &= (w^T w)^2 \neq
 \end{aligned}$$

Problem 5

⑤ loss of basic autoencoder: $\sum_{i=1}^d (g_i(x) - x_i)^2$.

let $h = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_d \end{bmatrix}$ represent the vector of hidden layer.

$$U_{ij} = W_{ij}^{(1)} = W_{ji}^{(2)}$$

$$g_i(x) = \sum_{n=1}^{\tilde{d}} W_{ni}^{(2)} h_n = \sum_{n=1}^{\tilde{d}} U_{in} h_n$$

$$= \sum_{n=1}^{\tilde{d}} U_{in} \tanh\left(\sum_{m=1}^d W_{mn}^{(1)} x_m\right)$$

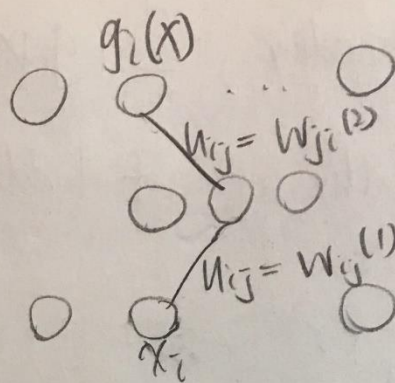
$$= \sum_{n=1}^{\tilde{d}} U_{in} \tanh\left(\sum_{m=1}^d U_{mn} x_m\right)$$

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

$$= \sum_{i=1}^d \left(\left(\sum_{n=1}^{\tilde{d}} U_{in} \tanh\left(\sum_{m=1}^d U_{mn} x_m\right) \right) - x_i \right)^2$$

Problem 6

(b)



let output of i th hidden neuron = h_i

$$E = \sum_{n=1}^d (g_n(x) - x_n)^2$$

$$\frac{\partial E}{\partial g_n(x)} = 2(g_n(x) - x_n) \quad \text{let } \delta_i = \frac{\partial(\text{output of } i\text{th hidden})}{\partial(\text{input of } i\text{th hidden})}$$

$$\text{if } n \neq i, \quad \frac{\partial g_n(x)}{\partial u_{ij}} = \frac{\partial(u_{ij} h_j)}{\partial u_{ij}} = h_j \delta_j x_i$$

$$\begin{aligned} \text{if } n = i \\ \frac{\partial g_n(x)}{\partial u_{ij}} &= \frac{\partial(u_{ij} h_j)}{\partial u_{ij}} = h_j + u_{ij} \frac{\partial h_j}{\partial u_{ij}} \\ &= h_j + u_{ij} \delta_j x_i \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial E}{\partial u_{ij}} &= \sum_{n=1}^d 2(g_n(x) - x_n) \cdot u_{ij} \delta_j x_i \\ &\quad + 2(g_i(x) - x_i) h_j \end{aligned}$$

$$\frac{\partial g_n(x)}{\partial w_{ij}^{(1)}} = \frac{\partial (w_{jn}^{(2)} h_j)}{\partial w_{ij}^{(1)}} = w_{jn}^{(2)} \delta_j x_i$$

$$= u_{nj} \delta_j x_i$$

$$\frac{\partial g_n(x)}{\partial w_{ji}^{(2)}} = \frac{\partial (w_{ji}^{(2)} h_j)}{\partial w_{ji}^{(2)}} = \begin{cases} h_j, & \text{if } n=i \\ 0, & \text{if } n \neq i \end{cases}$$

$$\therefore \frac{\partial E}{\partial w_{ij}^{(1)}} + \frac{\partial E}{\partial w_{ji}^{(2)}}$$

$$= \sum_{n=1}^d 2(g_n(x) - x_n) u_{nj} \delta_j x_i + 2(g_i(x) - x_i) h_j$$

$$= \frac{\partial E}{\partial u_{ij}} \Rightarrow \text{proved!}$$

Problem 7

⑦ hypothesis $g_{LIN}(x) = \text{sign}(w^T x + b)$

$w = x_+ - x_-$ and hyperplane $w^T x + b = 0$

passes through mid-point of x_+ & $x_- \rightarrow \frac{x_+ + x_-}{2}$

$$\therefore (x_+ - x_-)^T \left(\frac{x_+ + x_-}{2} \right) + b = 0 \quad \frac{1}{2} (\|x_+\|^2 - \|x_-\|^2) + b = 0$$

$$b = -\frac{1}{2} (\|x_+\|^2 - \|x_-\|^2)$$

$$\therefore g_{LIN}(x) = \text{sign} \left((x_+ - x_-)^T x - \frac{1}{2} (\|x_+\|^2 - \|x_-\|^2) \right)$$

Problem 8

⑧

$$g_{\text{RBFNET}} = \text{sign} \left(\beta_+ \exp(-\|x - \mu_+\|^2) + \beta_- \exp(-\|x - \mu_-\|^2) \right) \\ = \text{sign} \left(\exp(\|x - \mu_-\|^2 - \|x - \mu_+\|^2) + \frac{\beta_-}{\beta_+} \right).$$

$$\|x - \mu_-\|^2 - \|x - \mu_+\|^2 = (x - \mu_-)^T (x - \mu_-) - (x - \mu_+)^T (x - \mu_+) \\ = (x^T x - 2\mu_-^T x + \mu_-^T \mu_-) - (x^T x - 2\mu_+^T x + \mu_+^T \mu_+) \\ = 2(\mu_+ - \mu_-)^T x + (\mu_-^T \mu_- - \mu_+^T \mu_+)$$

$$g_{\text{RBFNET}} = \text{sign} \left(\|x - \mu_-\|^2 - \|x - \mu_+\|^2 - \ln\left(-\frac{\beta_-}{\beta_+}\right) \right) \\ = \text{sign} \left(2(\mu_+ - \mu_-)^T x + (\mu_-^T \mu_- - \mu_+^T \mu_+ - \ln\left(-\frac{\beta_-}{\beta_+}\right)) \right) \\ = \text{sign}(w^T x + b).$$

$$\text{Where } w = 2(\mu_+ - \mu_-), b = \mu_-^T \mu_- - \mu_+^T \mu_+ - \ln\left(-\frac{\beta_-}{\beta_+}\right)$$

Problem 9

⑨ $V_n = 1 \quad \forall n = 1, 2, \dots, N$ after initialization.

W_m is optimized via least square error over $\{V_n, r_{nm}\}$

$$\text{let error} = E = \sum_{n=1}^N (V_n W_m - r_{nm})^2 = \sum_{n=1}^N (W_m - r_{nm})^2$$

$$\frac{\partial E}{\partial W_m} = \sum_{n=1}^N 2(W_m - r_{nm})$$

$$\frac{\partial E}{\partial W_m} = 0 \quad \text{when} \quad \sum_{n=1}^N 2(W_m - r_{nm}) = 0$$

$$\Rightarrow W_m = \frac{1}{N} \sum_{n=1}^N r_{nm}$$

\Rightarrow average rating of the m -th movie

Problem 10

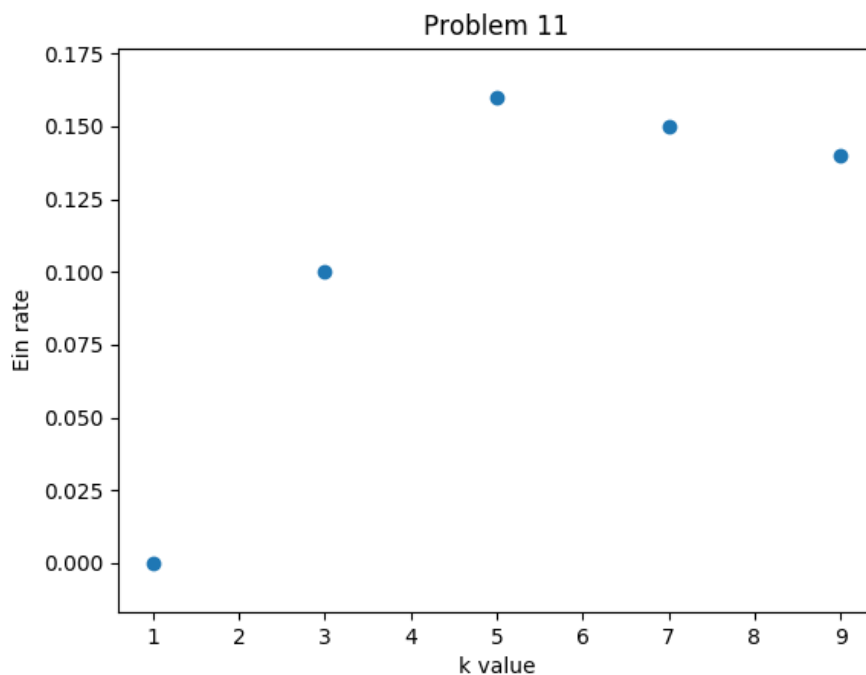
(19)

$$\begin{aligned} V_{n+1}^T W_m &= \left(\frac{1}{N} \sum_{n=1}^N V_n \right)^T W_m \\ &= \frac{1}{N} \sum_{n=1}^N V_n^T W_m = \frac{1}{N} \sum_{n=1}^N r_{nm} \end{aligned}$$

\equiv The average rating of the m th movie

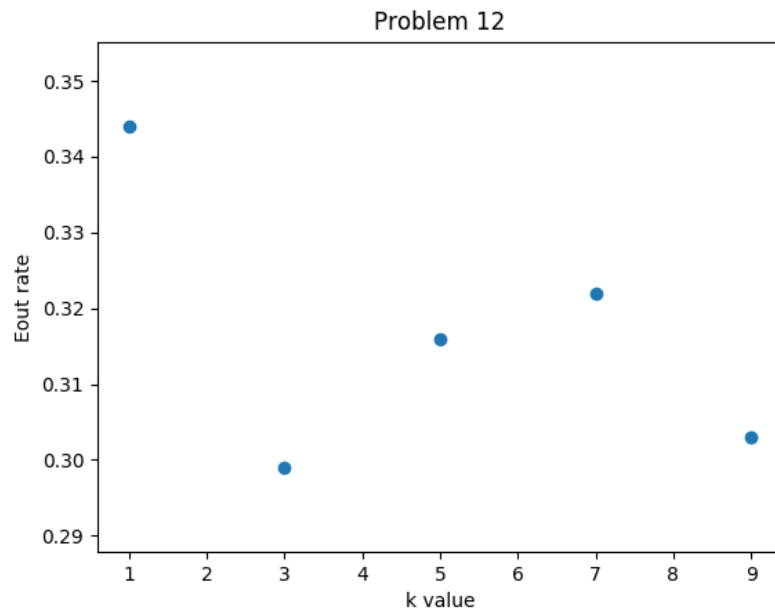
\therefore Maximum predicted score $V_{n+1}^T W_m$
 \Rightarrow the movie with largest average rating

Problem 11



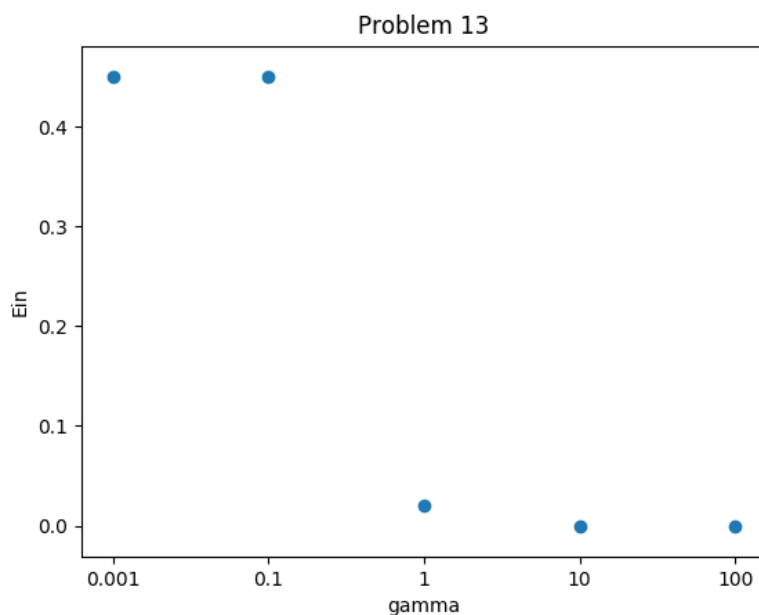
When $k = 1$, the prediction is obtained by the data itself. so the E_{in} is 0 as a result, and we can see that while k is increasing, the E_{in} rate is increasing, too.

Problem 12



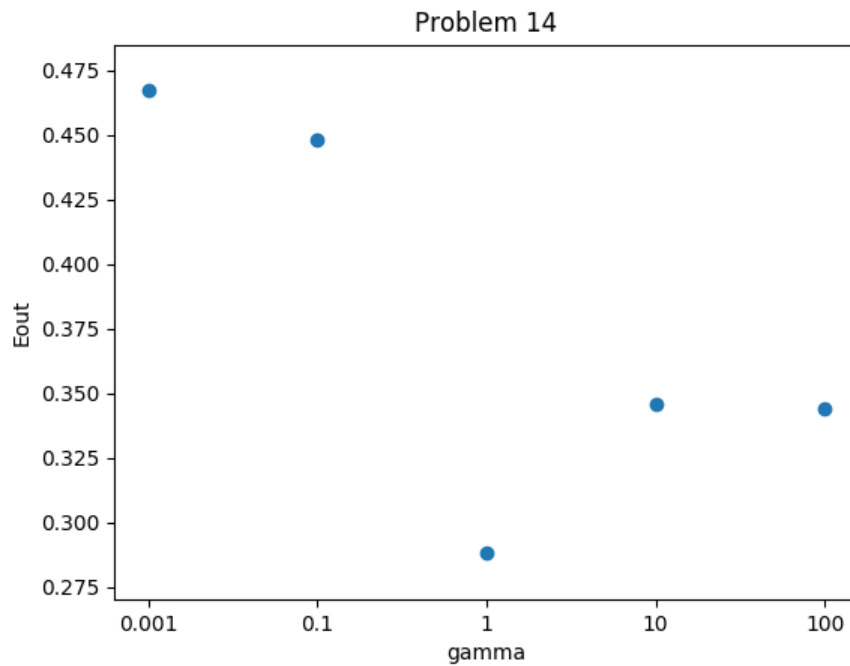
As we can see in this picture, when $k=1$, the prediction is obtained by the nearest data in the training data, and it's not sufficient; so when k gets bigger, the E out rate are smaller because the hypothesis consider more data around the testing data.

Problem 13



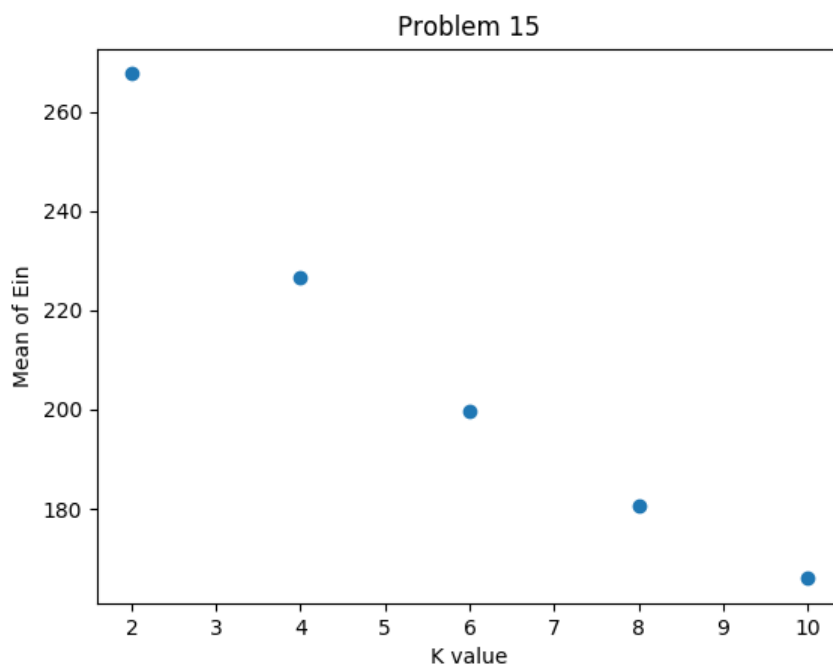
When gamma is small, the data far from the target can influence it largely, so the Error is high as a result, and when gamma gets bigger, the hypothesis tends to obtain prediction from the nearest data, so the error is low.

Problem 14



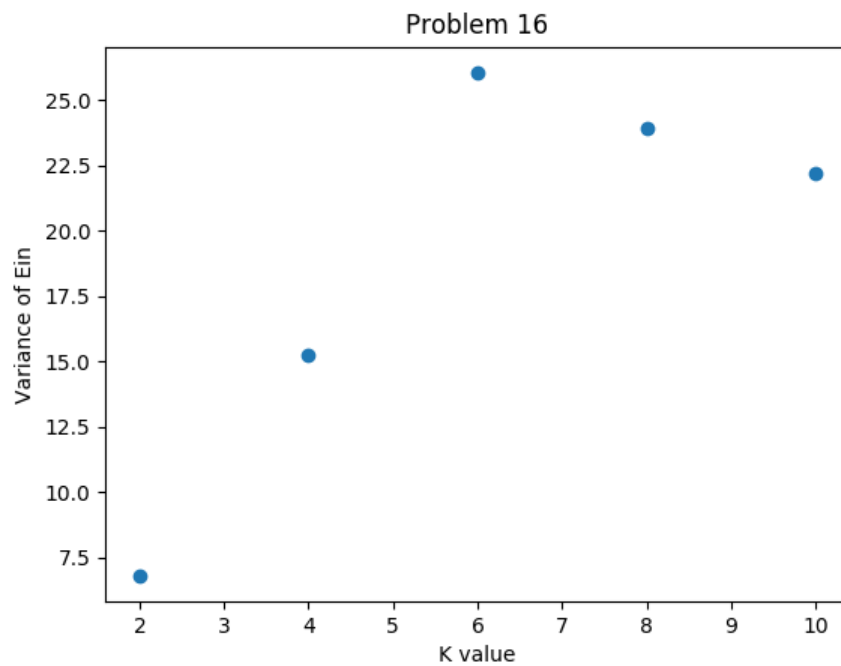
As we can see, when we chose a proper size for gamma (in this case = 1), can obtain the best performance.

Problem 15



While k value gets bigger, the E_{in} is decreasing, it's obvious because the average distance to nearest center is decreased.

Problem 16



The variance is increasing while k gets bigger, it's probably because the more centers can be learned, the more possible solution for convergence, which leads to larger variance.

Problem 17

(17)

consider extreme case: $N = 3\Delta \log_2 \Delta$

$$\frac{N^\Delta + 1}{2^N} = \frac{(3\Delta \log_2 \Delta)^\Delta + 1}{2^{3\Delta \log_2 \Delta}} = \frac{(3\Delta \log_2 \Delta)^\Delta + 1}{\Delta^{3\Delta}}$$

$$= \left(\frac{3\Delta \log_2 \Delta}{\Delta^3}\right)^\Delta + \frac{1}{\Delta^{3\Delta}} = \left(\frac{3 \log_2 \Delta}{\Delta^2}\right)^\Delta + \frac{1}{\Delta^{3\Delta}}$$

$$\text{let } f(\Delta) = \left(\frac{3 \log_2 \Delta}{\Delta^2}\right)^\Delta + \frac{1}{\Delta^{3\Delta}} = \left[\left(\frac{3}{\Delta}\right) \cdot \left(\frac{\log_2 \Delta}{\Delta}\right)\right]^\Delta + \frac{1}{\Delta^{3\Delta}}$$

$$\text{when } \Delta = 2, f(2) = \frac{9}{16} + \frac{1}{64} = \frac{37}{64} < 1.$$

$$\text{when } \Delta > 2, f(\Delta) < f(2) < 1. \quad (\because \frac{3}{\Delta}, \frac{\log_2 \Delta}{\Delta}, \frac{1}{\Delta^{3\Delta}} \text{ are decreasing})$$

$$\therefore \frac{N^\Delta + 1}{2^N} < 1 \text{ for all } \Delta \geq 2 \text{ when } N = 3\Delta \log_2 \Delta.$$

$$\therefore N^\Delta + 1 < 2^N \text{ for all } \Delta \geq 2 \text{ when } N = 3\Delta \log_2 \Delta.$$

$$\text{if } \Delta' < \Delta, N^{\Delta'} + 1 < N^\Delta + 1 < 2^N, \quad N = 3\Delta \log_2 \Delta > 3\Delta' \log_2 \Delta'$$

$$\therefore \text{for } \Delta \geq 2, \text{ if } N \geq 3\Delta \log_2 \Delta, \quad N^\Delta + 1 < 2^\Delta$$