

Mid Term Exam (Group MT-09) Machine Learning - 1 [CS/DS 864]

Darshan Bhat - MT2015038
Keerthan Pai K - MT2015053

October 19, 2016

Note: We have discussed with Group MT-03(Freeze Francis and Mohammed Haroon), but the representatons are different.

Q.1

The given description of H'_i s follows the structure of a feed-forward neural network with one hidden layer and one node in the output layer.

For given n input points, each H_i hypothesis class can produce $G_{H_i}(n)$ labellings. In the extreme case, the labellings produced are all distinct at the output of hidden layer. Now if H_1 can produce $G_{H_1}(n)$ different labellings of n input points, and H_2 can produce $G_{H_2}(n)$ different labellings independent of H_1 class ,then together we have $G_{H_1}(n) \times G_{H_2}(n)$ different possible labellings of n points at maximum.

With k hypothesis classes at the hidden layer, total number of distinct labellings is upper bounded by $\prod_{i=1}^k G_{H_i}(n)$. Now these labellings are fed to the output layer with H_0 hypothesis class. Each input can now be labelled with one of the $G_{H_0}(n)$ labellings. This implies a total of $G_{H_0}(n) \times \prod_{i=1}^k G_{H_i}(n)$ labellings are possible at the output layer.

That is,

$$G_H(n) \leq \prod_{i=0}^k G_{H_i}(n)$$

We need to prove

$$d_H \leq 2D \log_2 D$$

given that,

$$D > e \log_2 D$$

Proof:

$$\begin{aligned}
2^{d_H} &\leq \prod_{i=0}^k \left(\frac{ed_H}{d_{H_i}} \right)^{d_{H_i}} \\
&\leq \prod_{i=0}^k \left(\frac{ed_H}{2} \right)^{d_{H_i}} \\
&\leq \left(\frac{ed_H}{2} \right)^{\sum_{i=0}^k d_{H_i}} \\
2^{d_H} &\leq \left(\frac{ed_H}{2} \right)^D \\
d_H &\leq D \log_2 \left(\frac{ed_H}{2} \right)
\end{aligned}$$

This holds for all valid d_H values. We will prove our claim by contradiction.

Let $d_H > 2D \log_2 D$

i.e. $d_H = 2D \log_2 D + 1$, and we know that $D > e \log_2 D$

$$\begin{aligned}
\therefore 2D \log_2 D + 1 &\leq D \log_2 \left(\frac{e2D \log_2 D + e}{2} \right) \\
&\leq D \log_2 \left(\frac{2eD \log_2 D + 2eD \log_2 D}{2} \right) \\
&\leq D \log_2 \left(\frac{4eD \log_2 D}{2} \right) \\
&\leq D \log_2 (2eD \log_2 D) \\
\implies 2D \log_2 D &< D \log_2 (2eD \log_2 D) \\
2 \log_2 D &< \log_2 2 + \log_2 e + \log_2 D + \log_2 (\log_2 D) \\
\log_2 D &< \log_2 2 + \log_2 e + \log_2 (\log_2 D) \\
\implies \log_2 D &\leq \log_2 (e \log_2 D) \\
\implies D &\leq e \log_2 D
\end{aligned}$$

This contradicts with the given $D > e \log_2 D$.

Hence,

$$d_H \leq 2D \log_2 D$$

Q.2

2.a) We can prove that H matrix is idempotent.

$$\begin{aligned}
 H &= X (X^T X)^{-1} X^T \\
 H^2 &= X (X^T X)^{-1} X^T X (X^T X)^{-1} X^T \\
 H^2 &= X (X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T \\
 H^2 &= X (X^T X)^{-1} X^T \\
 H^2 &= H
 \end{aligned}$$

Now consider,

$$\begin{aligned}
 Hx &= \lambda x \\
 H^2x &= \lambda x \\
 HHx &= \lambda x \\
 H\lambda x &= \lambda x \\
 \lambda Hx &= \lambda x \\
 \lambda^2 x &= \lambda x \\
 (\lambda^2 - \lambda)x &= 0
 \end{aligned}$$

Solving for λ gives $\lambda = 0$ and $\lambda = 1$.

2.b)

The proof is by induction on the size n of the matrix A . The result is trivial for $n = 1$. Now let $n > 1$ and assume the result is true for any matrix of size $n - 1$.

Let λ be the eigenvalue of A , that is, $\det(\lambda I - A) = 0$, then $\lambda I - A$ must be non-invertible. This means that there exist a non-zero real vector u such that $Au = \lambda u$. We can always normalize u so that $u^T u = 1$. Thus, $\lambda = u^T A u$ is real. That is, the eigenvalues of a symmetric matrix are always real.

Now consider the eigenvalue λ_1 and an associated eigenvector u_1 . Using the Gram-Schmidt orthogonalization procedure, we can compute a $n \times (n - 1)$ matrix V_1 such that $[u_1, V_1]$ is orthogonal. By induction, we can write the $(n - 1) \times (n - 1)$ symmetric matrix $V_1^T A V_1$ as $Q_1 \Lambda_1 Q_1^T$, where Q_1 is a $(n - 1) \times (n - 1)$ matrix of eigenvectors, and $\Lambda_1 = \mathbf{diag}(\lambda_2, \dots, \lambda_n)$ are the $n - 1$ eigenvalues of $V_1^T A V_1$. Finally, we define the $n \times (n - 1)$ matrix $U_1 := V_1 Q_1$. By construction the matrix $U := [u_1, U_1]$ is orthogonal.

We have

$$U^T AU = \begin{pmatrix} u_1^T \\ U_1^T \end{pmatrix} A \begin{pmatrix} u_1 & U_1 \end{pmatrix} = \begin{pmatrix} u_1^T Au_1 & u_1^T AU_1 \\ U_1^T Au_1 & U_1^T AU_1 \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \Lambda_1 \end{pmatrix},$$

where we have exploited the fact that $U_1^T Au_1 = \lambda_1 U_1^T u_1 = 0$, and $U_1^T AU_1 = \Lambda_1$.

We have exhibited an orthogonal $n \times n$ matrix U such that $U^T AU$ is diagonal. This proves the theorem.

2.c)

The trace of a matrix is the sum of its diagonal entries. This has the property that $\text{Tr}(AB) = \text{Tr}(BA)$ for any two matrices of same size. $X^T X$ is a $(p+1) \times (p+1)$ matrix.

Now consider,

$$\begin{aligned} \text{Tr}(H) &= \text{Tr}(X(X^T X)^{-1} X^T) \\ &= \text{Tr}((X^T X)^{-1} X^T X) \\ &= \text{Tr}(I_{p+1}) \\ &= p+1 \\ \text{Tr}(H) &= \lambda_1 + \lambda_2 + \dots + \lambda_{p+1} \end{aligned}$$

Since $\lambda \in [0,1]$, we can conclude that number of eigen values which are 1's is $p+1$.

Q.3

To show that $X^t X$ is positive definite.
Consider,

$$v^t X^t X v = (Xv)^t (Xv) \quad (1)$$

Let $Xv = z$ with the dimension $n \times 1$.

$$(1) \Rightarrow z^t z = \|z\|^2 > 0$$

So, $X^t X > 0$ and thus it is positive semi definite.

We can model linear regression as as a system of linear equations. The vector equation is equivalent to a matrix equation of the form

$$X\mathbf{w} = \mathbf{y}$$

where X is an $n \times p$ matrix, \mathbf{w} is a column vector with p entries, and \mathbf{y} is a column vector with n entries.

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Solution for such a system exist only if \mathbf{y} is in the column space of \mathbf{X} . It may happend that there exist no solution. In such case we can find the best \mathbf{w} such that $\mathbf{X}\mathbf{w}$ is closest to \mathbf{b} by using least squares approximation.

We know that $\mathbf{X}\mathbf{w}$ is in the column space of \mathbf{X} and \mathbf{y} is not in the plane of $\mathbf{X}\mathbf{w}$. This can be achieved only when $\mathbf{X}\mathbf{w}$ is the projection of \mathbf{y}

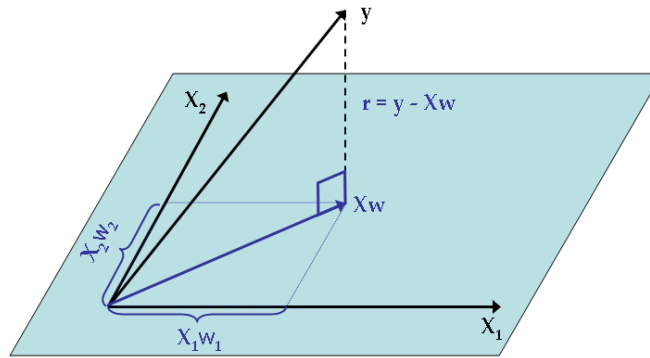


Figure 1:

$$\therefore \mathbf{X}\mathbf{w} \perp (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\begin{aligned}
\therefore (\mathbf{X}\mathbf{w})^t (y - \mathbf{X}\mathbf{w}) &= 0 \\
\mathbf{w}^t \mathbf{X}^t (y - \mathbf{X}\mathbf{w}) &= 0 \\
\mathbf{w}^t \mathbf{X}^t y - \mathbf{w}^t \mathbf{X}^t \mathbf{X} \mathbf{w} &= 0 \\
\mathbf{w}^t \mathbf{X}^t y &= \mathbf{w}^t \mathbf{X}^t \mathbf{X} \mathbf{w} \\
\mathbf{X}^t y &= \mathbf{X}^t \mathbf{X} \mathbf{w} \\
\mathbf{w} &= (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t y
\end{aligned}$$

Alternatively, the closed form solution can be derived using the positive semidefinite property of the matrix $\mathbf{X}^T \mathbf{X}$. Consider the loss function,

$$\begin{aligned}
S(\mathbf{w}) &= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \\
&= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\
&= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}
\end{aligned}$$

When $\mathbf{X}^T \mathbf{X}$ is positive definite, the quantity

$$S(\mathbf{w}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$$

can be written as

$$\langle \mathbf{w}, \mathbf{w} \rangle - 2\langle \mathbf{w}, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \rangle + \langle (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \rangle + C,$$

where C depends only on \mathbf{y} and \mathbf{X} , and $\langle \cdot, \cdot \rangle$ is the inner product defined by

$$\langle x, y \rangle = x^T (\mathbf{X}^T \mathbf{X}) y.$$

It follows that $S(\mathbf{w})$ is equal to

$$\langle \mathbf{w} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \mathbf{w} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \rangle + C$$

and therefore minimized exactly when

$$\mathbf{w} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = 0.$$

Hence the the best value of \mathbf{w} is,

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Q.4

4.a) y is of the form $y = g(x) + \epsilon_x$.

$$\text{Let } \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\text{So, } Y = g + \epsilon \tag{1}$$

$$\text{where } \mathbf{g} = \begin{bmatrix} g(x_1) \\ g(x_2) \\ \vdots \\ g(x_n) \end{bmatrix}, \quad \epsilon = \begin{bmatrix} \epsilon_{x_1} \\ \epsilon_{x_2} \\ \vdots \\ \epsilon_{x_n} \end{bmatrix}$$

$$\begin{aligned} E &= (Y^T - w^T X^T)(Y - Xw) \\ &= Y^T Y - Y^T Xw - w^T X^T Y + w^T X^T Xw \end{aligned}$$

To get optimal weights \hat{w} , $\frac{\delta E}{\delta w} = 0$

$$i.e. -X^T Y - X^T Y + 2X^T X \hat{w} = 0$$

$$X^T X \hat{w} = X^T Y$$

$$\therefore \hat{w} = (X^T X)^{-1} X^T Y$$

From equation (1),

$$\hat{w} = (X^T X)^{-1} X^T (g + \epsilon)$$

Estimated output $\hat{y} = X \hat{w}$

$$\hat{y} = X(X^T X)^{-1} X^T g + X(X^T X)^{-1} X^T \epsilon$$

From the given data, $g = Xw^*$

So,

$$\hat{Y} = X(X^T X)^{-1} X^T (Xw^*) + X(X^T X)^{-1} X^T \epsilon$$

$$\hat{Y} = Xw^* + \hat{H}\epsilon \tag{2}$$

From the given data,

$$y = x^T w^* + \epsilon$$

From (2) \hat{y} at a given point x will be

$$\hat{y} = x^T w^* + x^T (X^T X)^{-1} X^T \epsilon$$

$$y - \hat{y} = \epsilon - x^T (X^T X)^{-1} X^T \epsilon$$

4.b) Let

$$\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

We have,

$$\begin{aligned} x^T A x &= \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n \end{bmatrix} \\ &= a_{11}x_1^2 + a_{12}x_1x_2 + \cdots + a_{1n}x_1x_n + a_{21}x_2x_1 + a_{22}x_2^2 + \cdots + a_{2n}x_2x_n \\ &\quad + \cdots + a_{n1}x_nx_1 + \cdots + a_{nn}x_n^2 \end{aligned} \tag{1}$$

Also,

$$\begin{aligned} xx^T &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} = \begin{bmatrix} x_1^2 & x_1x_2 & \cdots & x_1x_n \\ x_2x_1 & x_2^2 & \cdots & x_2x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_nx_1 & x_nx_2 & \cdots & x_n^2 \end{bmatrix} \\ xx^T A &= \begin{bmatrix} x_1^2 & x_1x_2 & \cdots & x_1x_n \\ x_2x_1 & x_2^2 & \cdots & x_2x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_nx_1 & x_nx_2 & \cdots & x_n^2 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \\ \text{trace}(xx^T A) &= a_{11}x_1^2 + a_{21}x_1x_2 + \cdots + a_{n1}x_1x_n + a_{12}x_2x_1 + a_{22}x_2^2 + \cdots + a_{n2}x_2x_n \\ &\quad + \cdots + a_{n1}x_nx_1 + \cdots + a_{nn}x_n^2 \end{aligned} \tag{2}$$

Since A is symmetric, $a_{ij} = a_{ji}$

So, (1) = (2) that is,

$$x^T A x = \text{tr}(xx^T A)$$

Since Expectation $E(\cdot)$ and trace $\text{tr}(\cdot)$ both are linear operators,

$$\begin{aligned} \text{tr}[E[M]] &= \sum_i E[M_{ii}] \\ &= E\left[\sum_i M_{ii}\right] \\ &= E[\text{tr}[M]] \end{aligned}$$

4.c) x is a d -dimensional vector in R^d . We need to estimate the covariance matrix $E_D[xx^T]$ where expectation is over all possible datasets which is simply the expectation of xx^T in the domain of x . The matrix xx^T will be of the order $d \times d$.

$$xx^T = \begin{bmatrix} x_1^2 & x_1x_2 & \cdots & x_1x_d \\ x_2x_1 & x_2^2 & \cdots & x_2x_d \\ \vdots & \vdots & \ddots & \vdots \\ x_dx_1 & x_dx_2 & \cdots & x_d^2 \end{bmatrix}$$

$$E[xx^T] = \begin{bmatrix} E[x_1^2] & E[x_1x_2] & \cdots & E[x_1x_d] \\ E[x_2x_1] & E[x_2^2] & \cdots & E[x_2x_d] \\ \vdots & \vdots & \ddots & \vdots \\ E[x_dx_1] & E[x_dx_2] & \cdots & E[x_d^2] \end{bmatrix}$$

This matrix can be estimated by $X^T X$, where X is the $n \times d$ matrix. Each row of X is one sampled d -dimensional vector x_i . So we have n x_i samples arranged row wise in X . So,

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}$$

$$X^T = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}$$

$$X^T X_{d \times d} = \begin{bmatrix} (x_{11}^2 + x_{21}^2 + x_{31}^2 + \cdots) & (x_{11}x_{12} + x_{21}x_{22} + \cdots) & \cdots \\ (x_{12}x_{11} + x_{22}^2 + x_{32}x_{31} + \cdots) & (x_{12}^2 + x_{22}^2 + x_{32}^2 + \cdots) & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

The i, j^{th} entry in $X^T X$ will be

$$[X^T X]_{ij} = \sum_{k=1}^n x_{ki}x_{kj}$$

We see here that the i, j^{th} entry in $X^T X$ is filled with i, j^{th} entry in $[x_i x_j^T]$ matrix *for* $i = 1 \dots n$ i.e. we sample independent points from R^d space with some probability distribution (iid points). we calculate the underlying xx^T matrix for each of these points. Then we add every i, j^{th} from n xx^T matrices and copy to $[X^T X]_{ij}$ entry. Since the points are iid, clearly we are taking an average of $[xx^T]_{ij}$ with these samples.

$${}^i \begin{bmatrix} j \\ \cdot \\ a_{ij}^1 \end{bmatrix} \quad {}^i \begin{bmatrix} j \\ \cdot \\ a_{ij}^2 \end{bmatrix} \quad \cdots \quad {}^i \begin{bmatrix} j \\ \cdot \\ a_{ij}^n \end{bmatrix} \quad [X^T X] = \sum_{k=1}^n a_{ij}^k$$

from the law of large numbers $[X^T X]$ should converge to $nE_D[xx^T]$. From this we get the bound,

$$\begin{aligned} nE_D[xx^T] &= (1 + O(\sqrt[3]{n}))X^T X \\ \implies nE_D[xx^T] \cdot (X^T X)^{-1} &= (1 + O(\sqrt[3]{n}))I \end{aligned}$$

4.d) The true risk of linear regression with square error loss function is given by,

$$E_{D,\epsilon} [\|y - \hat{y}\|^2]$$

Consider

$$\begin{aligned} E_D [\|y - \hat{y}\|^2] &= E_D [(y - \hat{y})^T (y - \hat{y})] \\ &= E_D [(e - x^T (X^T X)^{-1} X^T \epsilon)^T (e - x^T (X^T X)^{-1} X^T \epsilon)] \\ &= E_D [(e^T - \epsilon^T X (X^T X)^{-1} x) (e - x^T (X^T X)^{-1} X^T \epsilon)] \\ &= E_D [e^T e - e^T x^T (X^T X)^{-1} X^T \epsilon - \epsilon^T X (X^T X)^{-1} x e + \epsilon^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon] \\ &= E_D \left[\text{tr} \begin{bmatrix} e^T e - e^T x^T (X^T X)^{-1} X^T \epsilon - \epsilon^T X (X^T X)^{-1} x e \\ + \epsilon^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon \end{bmatrix} \right] \quad (1) \end{aligned}$$

Consider,

$$\begin{aligned} \text{tr} [E_D [e^T x^T (X^T X)^{-1} X^T \epsilon]] &= \text{tr} [E_D [e^T x^T] (X^T X)^{-1} X^T \epsilon] \\ &= \text{tr} [E_D [e^T] E_D [x^T] (X^T X)^{-1} X^T \epsilon] \\ &= \text{tr} [0] \quad (\because E_D [e^T] = 0) \\ &= 0 \end{aligned}$$

Similarly,

$$\begin{aligned} \text{tr} [E_D [\epsilon^T X (X^T X)^{-1} x e]] &= \text{tr} [\epsilon^T X (X^T X)^{-1} E_D [x] E_D [e]] \\ &= 0 \end{aligned}$$

Therefore,

$$\begin{aligned} (1) &= E_D [\text{tr} [e^T e + \epsilon^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon]] \\ &= \sigma^2 + E_D \left[\text{tr} \left[\underbrace{\epsilon^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon}_{\text{tr} [xx^T (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}]} \right] \right] \\ &= \sigma^2 + E_D [\text{tr} [xx^T (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}]] \\ &= \sigma^2 + \text{tr} [E_D [xx^T] (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}] \end{aligned}$$

Now,

$$\begin{aligned}
E_\epsilon \left[E_D \left[\|y - \hat{y}\|^2 \right] \right] &= \sigma^2 + E_\epsilon \left[\text{tr} \left[E_D \left[xx^T \right] (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1} \right] \right] \\
&= \sigma^2 + \text{tr} \left[E_D \left[xx^T \right] (X^T X)^{-1} X^T E_\epsilon \left[\epsilon \epsilon^T \right] X (X^T X)^{-1} \right] \\
&= \sigma^2 + \text{tr} \left[E_D \left[xx^T \right] (X^T X)^{-1} X^T (\sigma^2 I) X (X^T X)^{-1} \right] \\
&= \sigma^2 + \frac{\sigma^2}{n} \text{tr} \left[n E_D \left[xx^T \right] \underline{(X^T X)^{-1} X^T X (X^T X)^{-1}} \right] \\
&= \sigma^2 + \frac{\sigma^2}{n} \text{tr} \left[n E_D \left[xx^T \right] (X^T X)^{-1} \right] \\
&= \sigma^2 + \frac{\sigma^2}{n} \left[\text{tr} \left(1 + o(\sqrt{n}) \right) I \right] \\
&= \sigma^2 + \frac{\sigma^2}{n} (d + 1 + o(\sqrt{n}) (d + 1)) \\
&= \sigma^2 \left(1 + \frac{d + 1}{n} + o\left(\frac{d + 1}{\sqrt{n}}\right) \right)
\end{aligned}$$

Q.5

Bounding $R(\hat{h})$ under $R_v(\hat{h})$ and intern $R(\hat{h}_l)$:

We are selecting \hat{h} from the set $H = \{\hat{h}_1, \hat{h}_2, \dots, \hat{h}_l, \dots, \hat{h}_k\}$ for which $R_v(\hat{h}_i)$ is the minimum. This is same as ERM principle for which the VC theorem bound holds.

So,

$$P\left(|R(\hat{h}) - R_v(\hat{h})| \geq \epsilon\right) \leq 4G_H(2n)e^{-2\epsilon^2(1-\alpha)n}$$

$G_H(2n)$ is bounded by $|H| = k$.

$$P\left(|R(\hat{h}) - R_v(\hat{h})| \geq \epsilon\right) \leq 4ke^{-2\epsilon^2(1-\alpha)n}$$

Let $4ke^{-2\epsilon^2(1-\alpha)n} = \frac{\delta}{2}$.

With probability of atleast $1 - \frac{\delta}{2}$,

$$R(\hat{h}) \leq R_v(\hat{h}) + \epsilon$$

Since $R_v(\hat{h}_l) \leq R_v(\hat{h})$,

$$R(\hat{h}) \leq R_v(\hat{h}_l) + \epsilon \tag{1}$$

Bounding $R_v(\hat{h}_l)$ in terms of $R(\hat{h}_l)$:

Since ERM principle holds for this case also,

$$P\left(|R_v(\hat{h}_l) - R(\hat{h}_l)| \geq \epsilon\right) \leq 4ke^{-2\epsilon^2(1-\alpha)n}$$

So with the probability of atleast $1 - \frac{\delta}{2}$,

$$R_v(\hat{h}_l) \leq R(\hat{h}_l) + \epsilon \tag{2}$$

Bounding $R(\hat{h}_l)$ in terms of $R_T(\hat{h}_l)$ in turn by $R_T(\hat{h}^*)$:

We selected \hat{h}_l from H_l class by ERM principle with training set risk as empirical risk. Here $G_{H_l}(2n)$ is bounded by c_l ,

$$\Rightarrow P\left(|R(\hat{h}_l) - R_T(\hat{h}_l)| \geq \epsilon'\right) \leq 4c_l e^{-2(\epsilon')^2 \alpha n}$$

Let $4c_l e^{-2(\epsilon')^2 \alpha n}$ be $\frac{\delta}{2}$.

Then with probability of atleast $\frac{\delta}{2}$,

$$R(\hat{h}_l) \leq R_T(\hat{h}_l) + \epsilon'$$

Since $R_T(\hat{h}_l) \leq R_T(h^*)$,

$$R(\hat{h}_l) \leq R_T(h^*) + \epsilon' \quad (3)$$

Bounding $R_T(h^*)$ using $R(h^*)$:

Again with V.C, theorem,

$$P(|R(h^*) - R_T(h^*)| \geq \epsilon') \leq 4c_l e^{-2(\epsilon')^2 \alpha n}$$

Then with probability of atleast $\frac{\delta}{2}$,

$$R_T(h^*) \leq R(h^*) + \epsilon' \quad (4)$$

From (1),(2),(3) and (4)

$$R(\hat{h}) \leq R(h^*) + 2\epsilon + 2\epsilon' \quad (5)$$

To calculate ϵ and ϵ' :

$$\begin{aligned} \frac{\delta}{2} &= 4k e^{-2\epsilon^2(1-\alpha)n} \\ \epsilon &= \sqrt{\frac{1}{2(1-\alpha)n} \ln \frac{8k}{\delta}} \end{aligned} \quad (6)$$

$$\begin{aligned} \frac{\delta}{2} &= 4c_l e^{-2(\epsilon')^2(\alpha)n} \\ \epsilon' &= \sqrt{\frac{1}{2(\alpha)n} \ln \frac{8c_l}{\delta}} \end{aligned} \quad (7)$$

Substituting (6) and (7) in (5),

$$R(\hat{h}) \leq R(h^*) + 2\sqrt{\frac{1}{2(1-\alpha)n} \ln \frac{8k}{\delta}} + 2\sqrt{\frac{1}{2(\alpha)n} \ln \frac{8c_l}{\delta}} \quad (8)$$

Both validation set equations $\{(1), (2)\}$ and training set equations $\{(3), (4)\}$ together hold with probability

$$(1 - \frac{\delta}{2}).(1 - \frac{\delta}{2}) = (1 - \frac{\delta}{2})^2 = 1 + \frac{\delta^2}{4} - \underline{\underline{2\frac{\delta}{2}}} \geq (1 - \delta)$$

With probability at least $(1 - \delta)$ equation (8) holds.