

# CSCC11-Assignment 1 Written Component

Masheil Mir

1005245600

MirMashe

Feb 4, 2021

1a) Consider LS basis function for single variable input/output

$$y = f(x) = w_0 + \sum_{k=1}^K w_k b_k(x)$$

Let  $\{(x_i, y_i)\}_{i=1}^N$  be the training data

$$y = f(x) = w_0 + w_1 b_1(x) + w_2 b_2(x) + \dots + w_K b_K(x)$$

Let  $\vec{x} = [x_1, \dots, x_N]^T$  be the N inputs

Let  $\vec{y} = [y_1, \dots, y_N]^T$  be the N outputs

Let  $\vec{w} = [w_0, w_1, \dots, w_N]^T$  be the weights

↳  $w_0$  is our bias term

$$B = \begin{bmatrix} 1 & b_1(x_1) & b_2(x_1) & \dots & b_K(x_1) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ i & b_1(x_N) & b_2(x_N) & \dots & b_K(x_N) \end{bmatrix}_{N \times (K+1)}$$

With least square objective  $f^2$ :

$$E(\vec{w}) = \sum_{i=1}^N (y_i - f(x_i))^2$$

$$= \sum_{i=1}^N [y_i - (w_0 + \sum_{k=1}^K w_k b_k(x_i))]^2$$

1b) Matrix form:

$$\begin{aligned} E(\vec{w}) &= \|\vec{y} - B\vec{w}\|_2^2 \\ &= (\vec{y} - B\vec{w})^T (\vec{y} - B\vec{w}) \\ &= (\vec{y}^T - \vec{w}^T B^T)(\vec{y} - B\vec{w}) \quad (1b) \text{ & } (1c) \text{ from matrix identities} \\ &= \vec{y}^T \vec{y} - \vec{w}^T B^T \vec{y} - \vec{y}^T B \vec{w} + \vec{w}^T B^T B \vec{w} \\ &= \vec{y}^T \vec{y} - 2\vec{w}^T B^T \vec{y} + \vec{w}^T B^T B \vec{w} \quad = 0 - 2B^T \vec{y} \end{aligned}$$

$$\begin{aligned} \nabla_{\vec{w}} E(\vec{w}) &= \frac{\partial E}{\partial \vec{w}} (\vec{y}^T \vec{y} - 2\vec{w}^T B^T \vec{y} + \vec{w}^T B^T B \vec{w}) \\ &= \underbrace{\frac{\partial}{\partial \vec{w}} (\vec{y}^T \vec{y})}_{=0} - 2 \underbrace{\frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T \vec{y})}_{\vec{w}^T B^T \vec{y}} + \underbrace{\frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T B \vec{w})}_{\vec{w}^T B^T B \vec{w}} \end{aligned}$$

Consider  $\frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T \vec{y})$

Let  $C = B^T \vec{y}$

$$\frac{\partial}{\partial \vec{w}} (\vec{w}^T C) \quad (6e)$$

$$= C = B^T \vec{y}$$

Consider  $\frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T B \vec{w})$

Let  $C = B^T B$

$$\frac{\partial}{\partial \vec{w}} (\vec{w}^T C \vec{w})$$

$$= (C + C^T) \vec{w} \quad (5b)$$

$$= (B^T B + (B^T B)^T) \vec{w}$$

$$= (B^T B + B^T B) \vec{w}$$

$$= 2B^T B \vec{w}$$

1b) Cont'd

$$\begin{aligned}\nabla_{\vec{w}} E(\vec{w}) &= \frac{\partial}{\partial \vec{w}} (\vec{y}^T \vec{y}) - 2 \frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T \vec{y}) + \frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T B \vec{w}) \\ &= -2 B^T \vec{y} + 2 B^T B \vec{w} \\ &= 2(B^T B \vec{w} - B^T \vec{y})\end{aligned}$$

1c) Optimal weight vector, denote this as  $\vec{w}^*$

$$\Rightarrow \frac{\partial E}{\partial \vec{w}} = 0$$

$$\Rightarrow 2(B^T B \vec{w} - B^T \vec{y}) = 0$$

$$\Rightarrow B^T B \vec{w} - B^T \vec{y} = 0$$

$$\Rightarrow B^T B \vec{w} = B^T \vec{y}$$

$$\Rightarrow \vec{w} = (B^T \vec{y}) (B^T B)^{-1}$$

$$\Rightarrow \vec{w}^* = (B^T \vec{y}) (B^T B)^{-1}$$

2a) In general, the following conditions will make the equation from (1c) not unique:

1) If the matrix  $B^T B$  is singular for the following equation:

$$\nabla_{\vec{w}} E(\vec{w}) = 0 \Rightarrow B^T B \vec{w} = B^T \vec{y}$$

i.e. Showing that the determinant is 0

2) When training data  $N < K$ ; not sufficient # of training data

$$y = f(x) = w_0 + \sum_{k=1}^K w_k b_k(x)$$

Let  $K=3$  and let our training data  $\{(x_i, y_i); i=0, 1, 2\}$ ,  $N=2$ ; Note  $N < K$

$$\text{Let } b_k(x) = x^k$$

$$f(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$\vec{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}; \quad \vec{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}; \quad B = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \end{bmatrix}; \quad \vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$B \vec{w} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} w_0 + w_1 x_1 + w_2 x_1^2 + w_3 x_1^3 \\ w_0 + w_1 x_2 + w_2 x_2^2 + w_3 x_2^3 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_0 + w_1 + w_2 + w_3 \end{bmatrix} \rightarrow \text{by plugging in } x_1=0, x_2=1 \text{ from training data}$$

$$B^T B \vec{w} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} w_0 \\ w_0 + w_1 + w_2 + w_3 \end{bmatrix} = \begin{bmatrix} 2w_0 + w_1 + w_2 + w_3 \\ w_0 + w_1 + w_2 + w_3 \\ w_0 + w_1 + w_2 + w_3 \\ w_0 + w_1 + w_2 + w_3 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$B^T \vec{y} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}; B^T B \vec{w} = B^T \vec{y} \Rightarrow \begin{bmatrix} 2 & 1 & 1 & 1 & | & 1 \\ 1 & 1 & 1 & 1 & | & 0 \\ 1 & 1 & 1 & 1 & | & 0 \\ 1 & 1 & 1 & 1 & | & 0 \end{bmatrix}$$

From the above it is shown that not having enough training data will make the normal eqs have multiple solutions since our system ends up being inconsistent (weights will have linear dependency).

$$2b) E(\vec{w}) = \underbrace{\|\vec{y} - B\vec{w}\|^2}_{\text{data term}} + \lambda \|\vec{w}\|^2 \quad \underbrace{\lambda \|\vec{w}\|^2}_{\text{Smoothness term}}$$

$$\begin{aligned} &= (\vec{y} - B\vec{w})^T (\vec{y} - B\vec{w}) + \lambda \vec{w}^T \vec{w} \\ &= (\vec{y} - B\vec{w})^T (\vec{y} - B\vec{w}) + \lambda \vec{w}^T \vec{w} \\ &= (\vec{y}^T - \vec{w}^T B^T) (\vec{y} - B\vec{w}) + \lambda \vec{w}^T \vec{w} \\ &= \vec{y}^T \vec{y} - \vec{w}^T B^T \vec{y} - \vec{y}^T B\vec{w} + \vec{w}^T B^T B\vec{w} + \lambda \vec{w}^T \vec{w} \\ &= \vec{y}^T \vec{y} - 2\vec{w}^T B^T \vec{y} + \vec{w}^T B^T B\vec{w} + \lambda \vec{w}^T \vec{w} \end{aligned}$$

$$\begin{aligned} \nabla_{\vec{w}} E(\vec{w}) &= \frac{\partial E}{\partial \vec{w}} (\vec{y}^T \vec{y} - 2\vec{w}^T B^T \vec{y} + \vec{w}^T B^T B\vec{w} + \lambda \vec{w}^T \vec{w}) \\ &= \frac{\partial}{\partial \vec{w}} (\vec{y}^T \vec{y}) - 2 \frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T \vec{y}) + \frac{\partial}{\partial \vec{w}} (\vec{w}^T B^T B\vec{w}) + \lambda \frac{\partial}{\partial \vec{w}} (\vec{w}^T \vec{w}) \\ &= 0 - 2B^T \vec{y} + 2B^T B\vec{w} + 2\lambda \vec{w} \\ &= 2(B^T B + \lambda I)\vec{w} - 2B^T \vec{y} \end{aligned}$$

$$\begin{aligned} \text{Set } \nabla_{\vec{w}} E(\vec{w}) &= 0 \Rightarrow 2(B^T B + \lambda I)\vec{w} - 2B^T \vec{y} = 0 \\ \Rightarrow & 2(B^T B + \lambda I)\vec{w} = 2B^T \vec{y} \\ \Rightarrow & (B^T B + \lambda I)\vec{w} = B^T \vec{y} \\ \Rightarrow & \vec{w}^* = (B^T B + \lambda I)^{-1} (B^T \vec{y}) \end{aligned}$$

Explanation:

This regularization helps ensure that  $\vec{w}^*$  has a unique value. The  $\lambda I$  term handles the issue of a Singular matrix. Note that  $\lambda > 0$  is given. The problem arises when we take  $B^T B$  being nearly singular. i.e.  $B^T B$  is positive semidefinite, it has eigenvalues  $\geq 0$ .  $\therefore B^T B + \lambda I$  ensures all eigenvalues strictly  $> 0$ , implying  $B^T B$  is invertible.

Thus this also means  $\vec{w}^*$  has only one unique solution.

2c)

$$\hat{B} = \begin{pmatrix} \vec{B} \\ \sqrt{\lambda} I_{k+1} \end{pmatrix} \in \mathbb{R}^{(N+k+1) \times (k+1)} ; \quad \hat{y} = \begin{pmatrix} \vec{y} \\ \vec{o}_{k+1} \end{pmatrix} \in \mathbb{R}^{(N+k+1)}$$

$$\hat{B} = \begin{bmatrix} 1 & b_1(x_1) & b_2(x_1) & \cdots & b_k(x_1) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & b_1(x_N) & b_2(x_N) & \cdots & b_k(x_N) \\ \sqrt{\lambda} & \sqrt{\lambda} & \ddots & & 0 \\ 0 & \ddots & & & \sqrt{\lambda} \end{bmatrix} ; \quad \hat{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ 0 \\ \vdots \\ 0 \end{bmatrix} ; \quad \vec{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix}$$

WTS:  $E(\vec{w}) = \|\hat{y} - \hat{B}\vec{w}\|_2^2$  is equivalent to  $E(\vec{w}) = \|\hat{y} - B\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2$

Consider  $E(w) = \|\hat{y} - \hat{B}\vec{w}\|_2^2$

$$= \left\| \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \begin{bmatrix} 1 & b_1(x_1) & b_2(x_1) & \cdots & b_k(x_1) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & b_1(x_N) & b_2(x_N) & \cdots & b_k(x_N) \\ \sqrt{\lambda} & \sqrt{\lambda} & \ddots & & 0 \\ 0 & \ddots & & & \sqrt{\lambda} \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix} \right\|_2^2$$

$$\begin{aligned} &= \sum_{i=1}^N y_i - w_0 - \left( \sum_{k=1}^K w_k b_k(x_i) \right)^2 + \sum_{k=0}^K (0 - \sqrt{\lambda} w_k)^2 \\ &= \|\hat{y} - B\vec{w}\|_2^2 + \sum_{k=0}^K (0 - \sqrt{\lambda} w_k)^2 \quad \text{by def of LS objective} \\ &= \|\hat{y} - B\vec{w}\|_2^2 + \lambda \sum_{k=0}^K (w_k)^2 \\ &= \|\hat{y} - B\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2, \quad \text{as wanted} \end{aligned}$$

$\therefore E(\vec{w}) = \|\hat{y} - \hat{B}\vec{w}\|_2^2$  is equivalent to  $E(\vec{w}) = \|\hat{y} - B\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2$ , as shown above, we can obtain the solution for regularized regression can be obtained by either method.

The constrain for  $\vec{w}$  can be found using  $\nabla_w E(w) = 0$  which was computed in part (2b) where we got the result of  $w^* = (B^T B + \lambda I)^{-1} (B^T \hat{y})$ .

3a) Let  $y = f(x) + \epsilon$ ,  $y \sim N(f(x), \sigma^2)$

Let  $\{(x_i, y_i)\}_{i=1}^N$  be the training data

Let  $f(x) = w_0 + w_1 b_1(x) + \dots + w_k b_k(x)$

Then, since  $y \sim N(f(x), \sigma^2)$ , we have

$$P(y | \vec{w}^\top \vec{b}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y - \vec{w}^\top \vec{b})^2}$$

$$\text{st } \vec{b} = [1 \quad b_1(x) \quad \dots \quad b_k(x)]$$

To compute ML objective:

$$P(y_{1:N} | \vec{w}^\top \vec{b}_1, \dots, \vec{w}^\top \vec{b}_N) \quad \text{where } \vec{b}_i = [1, b_1(x_i), \dots, b_k(x_i)]$$

$= \prod_{i=1}^N P(y_i | \vec{w}^\top \vec{b}_i)$  because each  $y_i$  is independent

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - \vec{w}^\top \vec{b}_i)^2} \quad \text{Since } y \sim N(f(x), \sigma^2)$$

$$= \left[ \frac{1}{\sqrt{2\pi}\sigma} \right]^N e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2}$$

$$= \left[ \frac{1}{2\pi\sigma^2} \right]^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2}$$

3b) The following defines an energy function to be minimized:

$$-\log(P(y_{1:N} | \vec{w}^\top \vec{b}_1, \dots, \vec{w}^\top \vec{b}_N))$$

$$= -\log \left( \left[ \frac{1}{2\pi\sigma^2} \right]^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2} \right)$$

$$= -\log \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} - \log e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2}$$

$$= -\left( \frac{N}{2} [\underbrace{\log(1)}_{=0} - \log(2\pi\sigma^2)] \right) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2$$

$$= \frac{N}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2$$

$$= C + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2 \quad \text{We can ignore } C \text{ since it does not affect the derivative \& multiply out } \frac{1}{2\sigma^2}$$

$$E(\vec{w}) = \sum_{i=1}^N (y_i - \vec{w}^\top \vec{b}_i)^2$$

$$= (\vec{y} - B\vec{w})^\top (\vec{y} - B\vec{w})$$

$$= \|\vec{y} - B\vec{w}\|_2^2, \quad B = [\vec{b}_1, \dots, \vec{b}_N]^\top$$

$$\text{Note: } N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

Recall LS objective from Q1:

$$E(\vec{w}) = \vec{y}^T \vec{y} - 2\vec{w}^T B^T \vec{y} + \vec{w}^T B^T B \vec{w} \rightarrow ①$$

Since we took the negative log likelihood, we find argmin to optimize  $\vec{w}^*$

Neg. log likelihood:  $E(\vec{w}) = \|y - B\vec{w}\|^2$

$$\begin{aligned} &= (\vec{y} - B\vec{w})^T (y - B\vec{w}) \\ &= (\vec{y}^T - \vec{w}^T B^T)(y - B\vec{w}) \\ &= \vec{y}^T \vec{y} - 2\vec{w}^T B^T \vec{y} + \vec{w}^T B^T B \vec{w} \rightarrow ② \end{aligned}$$

Since we can see the equation ① = ②, we can conclude that the negative log likelihood and the LS objective are the same.

3c)  $\vec{w} \sim N(0, \alpha^{-1} I)$ , Gaussian dist.

$$y \sim N(f(x), \alpha^{-1} I); \text{ Let } \vec{y} = \vec{w}^T \vec{b} + \varepsilon$$

Training data  $\{\vec{x}_i, y_i\}_{i=1}^N$

$$\begin{aligned} \text{Need } P(\vec{w} | \vec{b}_{1:N}, y_{1:N}) &= \frac{P(Y_{1:N} | \vec{b}_{1:N}, \vec{w}) P(\vec{w})}{P(Y_{1:N} | \vec{b}_{1:N})} \quad \text{Can ignore/disregard this term bc} \\ &\quad \text{doesn't depend on } \vec{w}. \\ &= P(Y_{1:N} | \vec{b}_{1:N}, \vec{w}) P(\vec{w}) \end{aligned}$$

$$① P(\vec{w}) = P(w_1, \dots, w_k)$$

$$= \prod_{i=1}^k P(w_i)$$

$$= \prod_{i=1}^k \frac{1}{\sqrt{2\pi\alpha^{-1}}} e^{(-\frac{\alpha}{2} \vec{w}_i^2)}$$

$$= \left( \frac{1}{\sqrt{2\pi\alpha^{-1}}} \right)^{N_2} e^{(-\frac{\alpha}{2} \vec{w}^T \vec{w})}$$

$$② \vec{y} = \vec{w}^T \vec{b} + \varepsilon$$

$$P(y_1, \dots, y_N | \vec{w}^T \vec{b}_1, \dots, \vec{w}^T \vec{b}_N) = \left[ \frac{1}{2\pi\sigma^2} \right]^{N_2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^T \vec{b}_i)^2} \quad (\text{Computed in part (3a)})$$

$$③ P(\vec{w} | \vec{b}_{1:N}, y_{1:N}) = P(Y_{1:N} | \vec{b}_{1:N}, \vec{w}) P(\vec{w})$$

$$= \left( \left[ \frac{1}{2\pi\sigma^2} \right]^{N_2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^T \vec{b}_i)^2} \right) \left( \left( \frac{1}{\sqrt{2\pi\alpha^{-1}}} \right)^{N_2} e^{(-\frac{\alpha}{2} \vec{w}^T \vec{w})} \right)$$

④ Negative log of MAP

$$\begin{aligned}
 -\log(P(\bar{w} | b_{1:N}, y_{1:N})) &= -\log \left[ \left( \left[ \frac{1}{2\pi\sigma^2} \right]^{N/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{w}^\top b_i)^2} \right) \left( \left( \frac{1}{\sqrt{2\pi}\alpha^{-1}} \right)^{N/2} e^{-\frac{\alpha}{2} \bar{w}^\top \bar{w}} \right) \right] \\
 &= \frac{N}{2} (-\log(1) - (-\log(2\pi\sigma^2)) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{w}^\top b_i)^2 + \frac{\alpha}{2} \bar{w}^\top \bar{w}) \\
 &= \frac{N}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{w}^\top b_i)^2 + \frac{\alpha}{2} \bar{w}^\top \bar{w} \\
 &= C + \frac{1}{2\sigma^2} \|y - B\bar{w}\|_2^2 + \frac{\alpha}{2} \|\bar{w}\|_2^2
 \end{aligned}$$

We can ignore  $C$ , multiply by  $2\sigma^2$

$$\Rightarrow E(\bar{w}) = \|y - B\bar{w}\|_2^2 + \sigma^2\alpha \|\bar{w}\|_2^2$$

3d) Looking at the negative log posterior compared to the LS objective above, we can conclude that they are the same, except in the negative log posterior we have the addition of the  $\sigma^2\alpha \|\bar{w}\|_2^2$  term which is the regularization term with  $\lambda = \sigma^2\alpha$ . This term aids with the issue of overfitting, and singularity of  $B^\top B$ .

That being said, we can say that the negative log posterior is almost identical to the regularized LS objective, but it has  $\lambda = \sigma^2\alpha$  in this case.

3e) If we assume  $\theta \sim U[0,1]$ , then  $P(\theta)$  is constant  $\forall \theta \in [0,1]$

$\therefore P(w)$  will not affect the min or max of the objective function and we can conclude that the MAP and ML objectives will be equivalent to the original LS objective.