

CSCC11 Assignment 1

Jefferson Li (jeffersonli.li@mail.utoronto.ca)

February 1, 2021

Question 1.

- a) Given training data denoted by

$$\{(x_i, y_i)_{i=1}^N\}$$

$$\begin{aligned} \text{let } \mathbf{w} &= [w_0, w_1, \dots, w_k]^T \\ \text{let } B &= \begin{bmatrix} 1 & b_1(x_1) & b_2(x_1) & \dots & b_k(x_1) \\ 1 & b_1(x_2) & b_2(x_2) & \dots & b_k(x_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & b_1(x_N) & b_2(x_N) & \dots & b_k(x_N) \end{bmatrix} \\ \text{let } \mathbf{y} &= [y_1, \dots, y_N]^T \end{aligned}$$

Then the LS objective can be formulated as such

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

$$E(\mathbf{w}) = \sum_{i=1}^N (y_i - f(x_i))^2 \quad \text{Definition of LS objective} \quad (2)$$

$$E(\mathbf{w}) = \sum_{i=1}^N (y_i - w_0 - \sum_{k=1}^K w_k b_k(x_i))^2 \quad (3)$$

$$E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 \quad \text{Replace with matrix norm representation} \quad (4)$$

$$E(\mathbf{w}) = (\mathbf{y} - B\mathbf{w})^T (\mathbf{y} - B\mathbf{w}) \quad \text{Property of Euclidean norm} \quad (5)$$

$$E(\mathbf{w}) = (\mathbf{y}^T - \mathbf{w}^T B^T) (\mathbf{y} - B\mathbf{w}) \quad \text{Properties of Transpose} \quad (6)$$

$$E(\mathbf{w}) = (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T B\mathbf{w} - \mathbf{w}^T B^T \mathbf{y} + \mathbf{w}^T B^T B\mathbf{w}) \quad (7)$$

$$E(\mathbf{w}) = (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T B^T \mathbf{y} + \mathbf{w}^T B^T B\mathbf{w}) \quad y^T Bw = w^T B^T y \text{ as they are both } [y_i B_{ij} x_j] \quad (8)$$

- b)

$$E(\mathbf{w}) = (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T B^T \mathbf{y} + \mathbf{w}^T B^T B\mathbf{w}) \quad (9)$$

$$\text{Gradient of } E = \frac{\partial E}{\partial \mathbf{w}} = (0 - 2B^T \mathbf{y} + (B^T B + B^T B)\mathbf{w}) \quad \text{From matrix identities 5a, 5b} \quad (10)$$

$$\frac{\partial E}{\partial \mathbf{w}} = -2B^T \mathbf{y} + 2B^T B\mathbf{w} \quad (11)$$

- c) The optimal weight vector \mathbf{w} is obtained by solving for \mathbf{w} when

$$\frac{\partial E}{\partial \mathbf{w}} = 0$$

$$-2B^T \mathbf{y} + 2B^T B\mathbf{w} = 0 \quad (12)$$

$$B^T B\mathbf{w} = B^T \mathbf{y} \quad (13)$$

$$(B^T B)^{-1} B^T B\mathbf{w} = (B^T B)^{-1} B^T \mathbf{y} \quad (14)$$

$$\mathbf{w}^* = (B^T B)^{-1} B^T \mathbf{y} = [w_0^*, w_1^*, \dots, w_k^*] \quad (15)$$

\implies The optimal weight vector \mathbf{w} is \mathbf{w}^*

Question 2.

- a) Q1(c) will not be unique if the columns of B are not linearly independent.

Proof. let B have dimensions m by n

Suppose Columns of B are not linearly independent

$$\implies \text{rank}(B) < \min(m, n) \quad (16)$$

$$\implies \text{rank}(B^T B) < \min(m, n) \quad \text{Since } \text{rank}(B) = \text{rank}(B^T B) \text{ by (2f)} \quad (17)$$

$$\implies B^T B \text{ is not invertible, and therefor has linearly dependent cols} \quad (18)$$

$$\implies \text{for } (B^T B)\mathbf{w} = B^T \mathbf{y}, \mathbf{w} \text{ has infinite solutions} \quad (19)$$

□

An example is if

$$K = 2, b_1 = x^2, b_2 = 2x^2$$

Note that column 1 and 2 of B are linearly dependent as column 2 of B is simply twice of column 1

Suppose the optimal weights are

$$w_0 = 1, w_1 = 2, w_2 = 3$$

So the optimal model is as follows

$$f(x) = 1 + 2b_1(x) + 3b_2(x) = 1 + 2x^2 + 6x^2 = 1 + 8x^2$$

This exact model can also be generated with the weights

$$w_0 = 1, w_1 = 8, w_2 = 0$$

as

$$f(x) = 1 + 8b_1(x) + 0b_2(x) = 1 + 8x^2 + 0x^2 = 1 + 8x^2$$

as we can see, if columns of B are not linearly independent, there is an infinite number of weights that correspond to the same model.

- b) Using the same $\mathbf{w}, B, \mathbf{y}$ definitions from Q1(a)

$$y = f(x) = w_0 + \sum_{k=1}^K w_k b_k(x) \quad (20)$$

$$E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \quad \text{Definition of Regularized LS objective} \quad (21)$$

$$E(\mathbf{w}) = (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T B^T \mathbf{y} + \mathbf{w}^T B^T B \mathbf{w} + \lambda \mathbf{w}^T \mathbf{w}) \quad \text{From the result of Q1(a)} \quad (22)$$

$$E(\mathbf{w}) = (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T B^T \mathbf{y} + \mathbf{w}^T [B^T B + \lambda I] \mathbf{w}) \quad (23)$$

$$\text{Gradient of } E = \frac{\partial E}{\partial \mathbf{w}} = (0 - 2B^T \mathbf{y} + ((B^T B + \lambda I) + (B^T B + \lambda I)^T) \mathbf{w}) \quad \text{From matrix identities 5a, 5b} \quad (24)$$

$$\frac{\partial E}{\partial \mathbf{w}} = (0 - 2B^T \mathbf{y} + (2B^T B + 2\lambda I) \mathbf{w}) \quad (25)$$

$$\frac{\partial E}{\partial \mathbf{w}} = 0 \implies 2B^T \mathbf{y} = (2B^T B + 2\lambda I) \mathbf{w} \quad (26)$$

$$\implies \mathbf{w}^* = (B^T B + \lambda I)^{-1} B^T \mathbf{y} \quad (27)$$

The regularization helps insure \mathbf{w}^* is a unique value.

Note that it is given that $\lambda > 0$

Since $B^T B$ is and is positive semidefinite, its eigenvalues are greater than or equal to zero (from the Linear Algebra Review and Reference pdf). Therefore $B^T B + \lambda I$ has all greater than zero eigenvalues. Which implies $B^T B + \lambda I$ is invertible.

Finally, this means

$$\mathbf{w}^* = (B^T B + \lambda I)^{-1} B^T \mathbf{y}$$

has one unique solution.

$$\bullet \text{ c) let } \hat{B} = \begin{bmatrix} 1 & b_1(x_1) & b_2(x_1) & \dots & b_k(x_1) \\ 1 & b_1(x_2) & b_2(x_2) & \dots & b_k(x_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & b_1(x_N) & b_2(x_N) & \dots & b_k(x_N) \\ \sqrt{\lambda} & & & & \\ & \sqrt{\lambda} & & & 0 \\ & & \ddots & & \\ & 0 & & \sqrt{\lambda} & \\ & & & & \sqrt{\lambda} \end{bmatrix} \quad \text{let } \hat{\mathbf{y}} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{let } \mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix}$$

Proof. Show $E(\mathbf{w}) = \|\hat{\mathbf{y}} - \hat{B}\mathbf{w}\|_2^2$ is equivalent to $E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$

$$E(\mathbf{w}) = \|\hat{y} - \hat{B}\mathbf{w}\|_2^2 \quad (28)$$

$$E(\mathbf{w}) = \left\| \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \frac{1}{\sqrt{\lambda}} \begin{bmatrix} 1 & b_1(x_1) & b_2(x_1) & \dots & b_k(x_1) \\ 1 & b_1(x_2) & b_2(x_2) & \dots & b_k(x_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & b_1(x_N) & b_2(x_N) & \dots & b_k(x_N) \\ \sqrt{\lambda} & & & 0 & \\ & \sqrt{\lambda} & & & \\ & & \ddots & & \\ 0 & & & \sqrt{\lambda} & \\ & & & & \sqrt{\lambda} \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix} \right\|_2^2 \quad (29)$$

$$E(\mathbf{w}) = \sum_{i=1}^N (y_i - w_0 - \sum_{k=1}^K w_k b_k(x_i))^2 + \sum_{k=0}^K (0 - \sqrt{\lambda} w_k)^2 \quad (30)$$

$$E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 + \sum_{k=0}^K (0 - \sqrt{\lambda} w_k)^2 \quad \text{By definition of LS Objective} \quad (31)$$

$$E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 + \lambda \sum_{k=0}^K (w_k)^2 \quad (32)$$

$$E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \quad (33)$$

as wanted. \square

Since I've shown that $E(\mathbf{w}) = \|\hat{y} - \hat{B}\mathbf{w}\|_2^2$ is equivalent to $E(\mathbf{w}) = \|\mathbf{y} - B\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$, the solution for regularized regression can be obtained using either function

Question 3.

- a)

let $\mathbf{w}^T \mathbf{b}(x_i) = f(x_i)$

$$E(w) = P(y_{1:N} | x_{1:N}, \mathbf{w}) \quad \text{Definition of Maximum Likelihood (ML) objective} \quad (34)$$

$$= \prod_{i=1}^N P(y_i | x_i, \mathbf{w}) \quad (35)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y_i - \mathbf{w}^T \mathbf{b}(x_i))^2}{2\sigma^2}\right\} \quad \text{Since } y \sim N(f(x), \sigma^2) \text{ and } \mathbf{w}^T \mathbf{b}(x_i) = f(x_i) \quad (36)$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^N \exp\left\{-\sum_{i=1}^N \frac{(y_i - \mathbf{w}^T \mathbf{b}(x_i))^2}{2\sigma^2}\right\} \quad (37)$$

$$= \left(\frac{1}{2\pi\sigma^2}\right)^{N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(x_i))^2\right\} \quad (38)$$

• b)

$$E(w) = \left(\frac{1}{2\pi\sigma^2}\right)^{N/2} \exp\left\{\frac{-1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2\right\} \quad (39)$$

$$-\ln(L(w)) = -\ln\left(\left(\frac{1}{2\pi\sigma^2}\right)^{N/2} \exp\left\{\frac{-1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2\right\}\right) \quad (40)$$

$$= -\left(\frac{N}{2} (\ln(1) - \ln(2\pi\sigma^2))\right) + \frac{-1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 \ln(e) \quad (41)$$

$$= \frac{N}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 \quad (42)$$

$$= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 \quad \text{Removing unnecessary constants} \quad (43)$$

$$= \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 \quad \text{Multiplying by the constant } \frac{1}{2\sigma^2} \quad (44)$$

$$= \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2 \quad \text{as } \mathbf{w}^T \mathbf{b}(x_i) = f(x_i) \quad (45)$$

$$(46)$$

Since we took the negative log likelihood, we look for the *argmin* when we optimize for \mathbf{w} .
Therefor this formulation is exactly the same as the LS Objective.

- c) let $1/K = P(y_{1:N}|x_{1:N})$

$$E(w) = P(\mathbf{w}|x_{1:N}, y_{1:N})$$

Definition of MAP

(47)

$$= \frac{P(y_{1:N}|x_{1:N}, \mathbf{w})p(\mathbf{w})}{P(y_{1:N}|x_{1:N})}$$

Bayes Rule

(48)

$$= K \cdot P(y_{1:N}|x_{1:N}, \mathbf{w})p(\mathbf{w})$$

Since $1/K = P(y_{1:N}|x_{1:N})$

(49)

$$= K \cdot P(y_{1:N}|x_{1:N}, \mathbf{w}) \left(\prod_{i=1}^K \frac{1}{\sqrt{2\pi\alpha^{-1}}} \exp\left\{ \frac{-1}{2\alpha^{-1}} w_i^2 \right\} \right)$$

Since $\mathbf{w} \sim N(0, \alpha^{-1}\mathbf{I})$

(50)

$$= K \cdot P(y_{1:N}|x_{1:N}, \mathbf{w}) \left(\left(\frac{1}{2\pi\alpha^{-1}} \right)^{N/2} \exp\left\{ \frac{-1}{2\alpha^{-1}} \sum_{i=1}^K w_i^2 \right\} \right)$$

(51)

$$-\ln(E(w)) = -\ln(K \cdot P(y_{1:N}|x_{1:N}, \mathbf{w}) \left(\left(\frac{1}{2\pi\alpha^{-1}} \right)^{N/2} \exp\left\{ \frac{-1}{2\alpha^{-1}} \sum_{i=1}^K w_i^2 \right\} \right))$$

negative ln both sides

(52)

$$= -\ln K - \ln P(y_{1:N}|x_{1:N}, \mathbf{w}) - \ln \left(\left(\frac{1}{2\pi\alpha^{-1}} \right)^{N/2} \right) - \ln \left(\exp\left\{ \frac{-1}{2\alpha^{-1}} \sum_{i=1}^K w_i^2 \right\} \right)$$

(53)

$$= -\ln P(y_{1:N}|x_{1:N}, \mathbf{w}) - \ln \left(\exp\left\{ \frac{-1}{2\alpha^{-1}} \sum_{i=1}^K w_i^2 \right\} \right)$$

Remove unnecessary constants

(54)

$$= -\ln P(y_{1:N}|x_{1:N}, \mathbf{w}) + \frac{1}{2\alpha^{-1}} \sum_{i=1}^K w_i^2$$

(55)

$$= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 + \frac{1}{2\alpha^{-1}} \sum_{i=1}^K w_i^2$$

 $-\ln P(y_{1:N}|x_{1:N}, \mathbf{w})$ from part b

(56)

$$= \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 + \frac{\sigma^2}{\alpha^{-1}} \sum_{i=1}^K w_i^2$$

multiply by constant $2\sigma^2$

(57)

$$= \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{b}(\mathbf{x}_i))^2 + (\sigma^2 \alpha) \|\mathbf{w}^T \mathbf{w}\|_2^2$$

(58)

- d) Since we took the negative log MAP, we look for the *argmin* when we optimize for \mathbf{w} .
Therefor this formulation is exactly the same as the Regularized LS Objective (ridge regression) where $\lambda = \sigma^2 \alpha$.
- e) Since for a uniform distribution, $P(\mathbf{w})$ is constant for all $\mathbf{w} \in [l, u]$
Therefor, $P(\mathbf{w})$ has no effect on the minimum or maximum of the objective function. so the MAP and ML objectives would both be equivalent to the LS objective.