

Theoretical (50)

1 (6 points)

Part1: (c) A high magnitude suggests that the feature is important. However, it may be the case that another feature is highly correlated with this feature and it's coefficient also has a high magnitude with the opposite sign, in effect cancelling out the effect of the former. Thus, we cannot really remark on the importance of a feature just because it's coefficient has a relatively large magnitude.

Part2:

- True , we can not reduce the bias and cover our model's limits by increasing the training data.
- False , the obvious counterexample that can be mentioned is overfitting on training data.
- False , increasing model complexity could be useful for this case , but not always.
- False

2 (19 points)

2.a:

$$\begin{aligned}
 J(w) &= \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2 \\
 \nabla_w J(w) &= -2 \sum_{i=1}^n (y^{(i)} - w^T x^{(i)}) (x^{(i)})^T = 0 \\
 \sum_{i=1}^n y^{(i)} (x^{(i)})^T &= w^T \sum_{i=1}^n x^{(i)} (x^{(i)})^T \\
 x^T y &= (x^T x) w \\
 w^* &= (x^T x)^{-1} x^T y
 \end{aligned}$$

2.b:

$$\begin{aligned}
 J(w) &= \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2 + \lambda w^T w \\
 \nabla_w J(w) &= -2 \sum_{i=1}^n (y^{(i)} - w^T x^{(i)}) (x^{(i)})^T + 2\lambda \sum_{j=1}^M w_j = 0 \\
 \sum_{i=1}^n y^{(i)} (x^{(i)})^T &= w^T \sum_{i=1}^n x^{(i)} (x^{(i)})^T + \lambda \sum_{j=1}^M w_j \\
 x^T y &= (x^T x) w + \lambda I w \\
 x^T y &= (x^T x + \lambda I) w \\
 w^* &= (x^T x + \lambda I)^{-1} x^T y
 \end{aligned}$$

2.c:

First, we should prove that if $\Sigma X = XF$ then $w_{new}^* = w^*$

$$\begin{aligned}
 \Sigma X = XF &\Rightarrow \Sigma = XF X^{-1} \Rightarrow \Sigma^{-1} = XF^{-1} X^{-1} \\
 &\Rightarrow \Sigma^{-T} = X^{-T} F^{-T} X^T \Rightarrow \Sigma^{-1} = X^{-T} F^{-T} X^T \\
 w_{new}^* &= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y = (X^T X^{-T} F^{-T} X^T X)^{-1} X^T X^{-T} F^{-T} X^T y \\
 &= (X^T X)^{-1} F^T F^{-T} X^T y = (X^T X)^{-1} X^T y = w^*
 \end{aligned}$$

The first side of the term was proved.

Then, we should prove that if $w_{new}^* = w^*$ then $\Sigma X = XF$

$$\begin{aligned}
 w^* &= (X^T X)^{-1} X^T y = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y \Rightarrow (X^T X)^{-1} X^T = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \\
 X^T \Sigma &= (X^T X) (X^T \Sigma^{-1} X)^{-1} X^T \Rightarrow \Sigma X = X [(X^T X) (X^T \Sigma^{-1} X)^{-1}]^T
 \end{aligned}$$

now name $[(X^T X) (X^T \Sigma^{-1} X)^{-1}]^T$ as F then we conclude the other side so this estimator is equal to ordinary least-square estimator w_{opt}^* , if and only if there exist a nonsingular matrix F , such that $\Sigma X = XF$.

4 (15 points)

4.a

In this particular case, it looks like our data matrix is equal to x_j . So according to the linear regression relation we have equality below:

$$w_j = (x_j^T x_j)^{-1} x_j^T y = \frac{x_j^T y}{x_j^T x_j}$$

4.b

Note that the columns are orthogonal therefore their internal multiplication is zero. As a result, the value of $X^T X$ will be diagonal. More precisely:

$$X^T X = \text{diag}(x_1^T x_1, \dots, x_m^T x_m) \Rightarrow (X^T X)^{-1} = \text{diag}((x_1^T x_1)^{-1}, \dots, (x_m^T x_m)^{-1})$$

now we have:

$$\begin{aligned} w &= (X^T X)^{-1} X^T y = \text{diag}((x_1^T x_1)^{-1}, \dots, (x_m^T x_m)^{-1}) X^T y \Rightarrow \\ w_j &= (\text{diag}((x_1^T x_1)^{-1}, \dots, (x_m^T x_m)^{-1}) X^T y)_j = (\text{diag}((x_1^T x_1)^{-1}, \dots, (x_m^T x_m)^{-1}))_j (X^T y)_j \\ &= (x_j^T x_j)^{-1} (X^T y)_j = \frac{(X^T y)_j}{x_j^T x_j} = \frac{(X^T)_j y}{x_j^T x_j} = \frac{x_j^T y}{x_j^T x_j} \end{aligned}$$

Combining the recent equality with the previous part, results to concluding that optimal parameters from training the regressor on all features is the same as the optimal parameters resulting from training on each feature independently".

4.c

Note that in the matrix mode our data is as follows:

$$X = \begin{bmatrix} 1 \\ 1 \\ x_j \\ \cdot \\ \cdot \\ 1 \end{bmatrix}$$

That x_j is a column. For the internal multiplication of $(x_j)(1, \dots, 1)$, which is equal to the sum of the x_j elements. We use the $\text{sum}(x_j)$ symbol. Now with this in mind we have:

$$X^T X = \begin{bmatrix} x_j^T x_j & \text{sum}(x_j) \\ \text{sum}(x_j) & n \end{bmatrix} \Rightarrow (X^T X)^{-1} = \frac{1}{n \|x_j\|^2 - \text{sum}(x_j)^2} \begin{bmatrix} n & -\text{sum}(x_j) \\ -\text{sum}(x_j) & x_j^T x_j \end{bmatrix}$$

Note that:

$$X^T y = \left(\begin{bmatrix} 1 \\ 1 \\ x_j \\ \cdot \\ \cdot \\ 1 \end{bmatrix} \right)^T y = \begin{bmatrix} x_j^T y \\ \text{sum}(y) \end{bmatrix}$$

Then:

$$[w_j, w_0] = w = (X^T X)^{-1} X^T y = \frac{1}{n \|x_j\|^2 - \text{sum}(x_j)^2} \begin{bmatrix} n & -\text{sum}(x_j) \\ -\text{sum}(x_j) & x_j^T x_j \end{bmatrix} \begin{bmatrix} x_j^T y \\ \text{sum}(y) \end{bmatrix}$$

$$w_j = \frac{n x_j^T y - \text{sum}(x_j) \text{sum}(y)}{n \|x_j\|^2 - \text{sum}(x_j)^2} = \frac{\frac{x_j^T y}{n} - \frac{\text{sum}(x_j)}{n} \frac{\text{sum}(y)}{n}}{\frac{\|x_j\|^2}{n} - \left(\frac{\text{sum}(x_j)}{n}\right)^2} = \frac{\mathbf{E}[x_j y] - \mathbf{E}[x_j] \mathbf{E}[y]}{\mathbf{E}[x_j^2] - \mathbf{E}[x_j]^2} = \frac{\text{cov}(x_j, y)}{\text{var}(x_j)}$$

$$\begin{aligned} w_0 &= \frac{\text{sum}(y) \|x_j\|^2 - \text{sum}(x_j) (x_j^T y)}{n \|x_j\|^2 - \text{sum}(x_j)^2} = \frac{\text{sum}(y) \|x_j\|^2 - \text{sum}(x_j) (x_j^T y)}{n \|x_j\|^2 - \text{sum}(x_j)^2} = \frac{n^2 \mathbf{E}[y] \mathbf{E}[x_j^2] - n^2 \mathbf{E}[x_j] \mathbf{E}[x_j y]}{n^2 \text{var}(x_j)} \\ &= \frac{\mathbf{E}[y] \mathbf{E}[x_j^2] - \mathbf{E}[x_j] \mathbf{E}[x_j y]}{\text{var}(x_j)} = \frac{\mathbf{E}[y] (\mathbf{E}[x_j^2] - \mathbf{E}[x_j]^2) + \mathbf{E}[y] \mathbf{E}[x_j]^2 - \mathbf{E}[x_j] \mathbf{E}[x_j y]}{\text{var}(x_j)} = \\ &= \mathbf{E}[y] + \frac{\mathbf{E}[y] \mathbf{E}[x_j]^2 - \mathbf{E}[x_j] \mathbf{E}[x_j y]}{\text{var}(x_j)} = \mathbf{E}[y] + \mathbf{E}[x_j] \frac{\mathbf{E}[y] \mathbf{E}[x_j] - \mathbf{E}[x_j y]}{\text{var}(x_j)} = \mathbf{E}[y] - w_j \mathbf{E}[x_j] \end{aligned}$$