

Artificial Intelligence: HW 2

Jeong Min Lee

November 6, 2023

1 Linear Regression

1.a

Let f be the target function. I'll use the superscript with parenthesis to describe the n-th sample vector.

$$f(\boldsymbol{\omega}, \overline{\mathbf{x}}^{(n)}) = \sum_n \frac{1}{2} \left(t^{(n)} - \boldsymbol{\omega}^T \overline{\mathbf{x}}^{(n)} \right)^2 \quad (1)$$

To minimize f , differentiate it by $\boldsymbol{\omega}$ and find the $\boldsymbol{\omega}_0$ which makes the derivative zero. To make expression simple, I used the Einstein notation.

$$\begin{aligned} \frac{\partial f}{\partial \omega_j} &= -(t^{(n)} - \boldsymbol{\omega}^T \overline{\mathbf{x}}^{(n)}) \cdot \frac{\partial}{\partial \omega_j} \boldsymbol{\omega}^T \overline{\mathbf{x}}^{(n)} \\ &= -(t^{(n)} - \boldsymbol{\omega}^T \overline{\mathbf{x}}^{(n)}) x_j^{(n)} \\ &= 0 \end{aligned}$$

By enumerating the $\frac{\partial f}{\partial \omega_j}$ horizontally, one can get $\frac{\partial f}{\partial \boldsymbol{\omega}}$.

$$\sum_n t^{(n)} \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_M^{(n)} \end{pmatrix}^T = \sum_n \begin{pmatrix} \boldsymbol{\omega}^T \overline{\mathbf{x}}^{(n)} x_1^{(n)} \\ \vdots \\ \boldsymbol{\omega}^T \overline{\mathbf{x}}^{(n)} x_M^{(n)} \end{pmatrix}^T \quad (2)$$

The left hand side is simply $\sum_n t^{(n)} \overline{\mathbf{x}}^{(n)T}$. From the linearity of vector summation rule, the right hand side is $\left(\left(\sum_n \overline{\mathbf{x}}^{(n)} \cdot \overline{\mathbf{x}}^{(n)T} \right) \boldsymbol{\omega} \right)^T$. By taking transpose to both sides, one can get the following equation.

$$\left[\sum_n \overline{\mathbf{x}}^{(n)} \cdot \overline{\mathbf{x}}^{(n)T} \right] \boldsymbol{\omega} = \sum_n t^{(n)} \overline{\mathbf{x}}^{(n)} \quad (3)$$

Therefore, $\mathbf{A} = \sum_n \overline{\mathbf{x}}^{(n)} \cdot \overline{\mathbf{x}}^{(n)T}$ and $\mathbf{b} = \sum_n t^{(n)} \overline{\mathbf{x}}^{(n)}$

1.b

$\overline{\mathbf{x}}^{(1)} = (1, 0)^T, t^{(1)} = 1. \overline{\mathbf{x}}^{(2)} = (1, \epsilon)^T, t^{(2)} = 1. \mathbf{A} = \overline{\mathbf{x}}^{(1)} \cdot \overline{\mathbf{x}}^{(1)T} + \overline{\mathbf{x}}^{(1)} \cdot \overline{\mathbf{x}}^{(2)T} = \begin{pmatrix} 2 & \epsilon \\ \epsilon & \epsilon^2 \end{pmatrix}$
 $\mathbf{b} = \overline{\mathbf{x}}^{(1)} + \overline{\mathbf{x}}^{(2)} = (2, \epsilon)^T$. Since \mathbf{A} is invertible (determinant is nonzero.),

$$\boldsymbol{\omega} = \mathbf{A}^{-1} \mathbf{b} \quad (4)$$

$$= \frac{1}{\epsilon^2} \begin{pmatrix} \epsilon^2 & -\epsilon \\ -\epsilon & 2 \end{pmatrix} \begin{pmatrix} 2 \\ \epsilon \end{pmatrix} \quad (5)$$

$$= \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad (6)$$

1.c

\mathbf{A} is same to the above one. $\mathbf{b} = (1 + \epsilon) \cdot \overline{\mathbf{x}^{(1)}} + \overline{\mathbf{x}^{(2)}} = (2 + \epsilon, \epsilon)^T$

$$\boldsymbol{\omega} = \mathbf{A}^{-1}\mathbf{b} \quad (7)$$

$$= \frac{1}{\epsilon^2} \begin{pmatrix} \epsilon^2 & -\epsilon \\ -\epsilon & 2 \end{pmatrix} \begin{pmatrix} 2 + \epsilon \\ \epsilon \end{pmatrix} \quad (8)$$

$$= \begin{pmatrix} 1 + \epsilon \\ -1 \end{pmatrix} \quad (9)$$

1.d

$\boldsymbol{\omega}_b = (1, 0)^T, \boldsymbol{\omega}_c = (1.1, -1)^T$. The difference of $\Delta\boldsymbol{\omega} = \boldsymbol{\omega}_c - \boldsymbol{\omega}_b = (\epsilon, -1)^T = (0.1, -1)^T$

2 Linear Regression with Regularization

2.a

Claim 1 : \mathbf{A} is positive semi-definite.

proof

\mathbf{A} is trivially symmetry matrix. $\forall \mathbf{v} \in \mathbb{R}^n, \mathbf{v}^T \mathbf{A} \mathbf{v} = \sum_n \mathbf{v}^T \overline{\mathbf{x}^{(n)}} \cdot \overline{\mathbf{x}^{(n)}}^T \mathbf{v} = \sum_n \|\mathbf{v}^T \overline{\mathbf{x}^{(n)}}\|^2 \geq 0. \blacksquare$

Claim 2 : $\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \iff \mathbf{A}^{-1}\mathbf{x} = \lambda^{-1}\mathbf{x}$ where $\lambda \neq 0$ and \mathbf{A} is invertible.

proof

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \iff \mathbf{A}^{-1}\mathbf{A}\mathbf{x} = \lambda\mathbf{A}^{-1}\mathbf{x} \iff \lambda^{-1}\mathbf{x} = \mathbf{A}^{-1}\mathbf{x}. \blacksquare$$

Let $S(\mathbf{A})$ be the set of all eigenvalues of \mathbf{A} .

$$S(\mathbf{A}) \equiv \{\lambda_i | \text{for some } \mathbf{x} \in \mathbb{R}, \mathbf{A}\mathbf{x} = \lambda_i\mathbf{x}\} \quad (10)$$

$\forall \tilde{\lambda} \in S(\mathbf{A} + \lambda\mathbf{I})$ s.t. $(\mathbf{A} + \lambda\mathbf{I})\mathbf{x} = \tilde{\lambda}\mathbf{x}$.

By multiplying $\mathbf{x}^T, \mathbf{x}^T(\mathbf{A} + \lambda\mathbf{I})\mathbf{x} = \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda = \tilde{\lambda} \geq \lambda$. ($\because \mathbf{A}$ is positive semi-definite.)

This implies that $\min(S(\mathbf{A} + \lambda\mathbf{I})) \geq \lambda$. Equivalently, due to the **Claim 2**, this also means that $\max(S((\mathbf{A} + \lambda\mathbf{I})^{-1})) \leq \lambda^{-1}$. By noticing that $\max(S((\mathbf{A} + \lambda\mathbf{I})^{-1})) = \rho((\mathbf{A} + \lambda\mathbf{I})^{-1})$, the proof is done. Note that for the equality, $\mathbf{A}\mathbf{x} = \mathbf{0}$ must have nontrivial solution.

2.b

For both problems, the $\mathbf{A} + \lambda\mathbf{I}$ is following.

$$\mathbf{A} + \lambda\mathbf{I} = \begin{pmatrix} 2 + \lambda & \epsilon \\ \epsilon & \epsilon^2 + \lambda \end{pmatrix} \quad (11)$$

Since (11) is invertible, one can get $\boldsymbol{\omega}_b, \boldsymbol{\omega}_c$.

$$\begin{aligned} \boldsymbol{\omega}_b &= \frac{1}{(1 + \lambda)\epsilon^2 + \lambda(\lambda + 2)} \cdot \begin{pmatrix} \epsilon^2 + \lambda & -\epsilon \\ -\epsilon & 2 + \lambda \end{pmatrix} \cdot \begin{pmatrix} 2 \\ \epsilon \end{pmatrix} \\ &= \frac{1}{(1 + \lambda)\epsilon^2 + \lambda(\lambda + 2)} \cdot \begin{pmatrix} \epsilon^2 + 2\lambda \\ \epsilon\lambda \end{pmatrix} \\ &= \begin{pmatrix} 0.973 \\ 0.044 \end{pmatrix} \end{aligned}$$

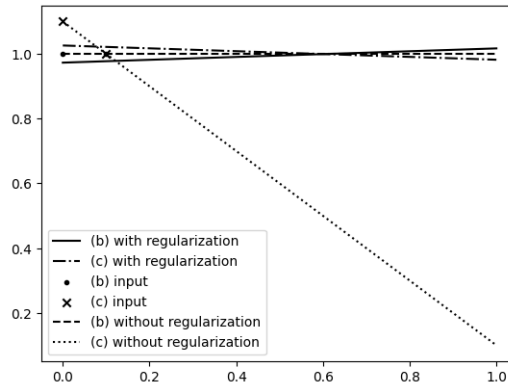
$$\begin{aligned}
\omega_c &= \frac{1}{(1+\lambda)\epsilon^2 + \lambda(\lambda+2)} \cdot \begin{pmatrix} \epsilon^2 + \lambda & -\epsilon \\ -\epsilon & 2 + \lambda \end{pmatrix} \cdot \begin{pmatrix} 2 + \epsilon \\ \epsilon \end{pmatrix} \\
&= \frac{1}{(1+\lambda)\epsilon^2 + \lambda(\lambda+2)} \cdot \begin{pmatrix} \epsilon^3 + \epsilon^2 + \lambda\epsilon + 2\lambda \\ -\epsilon^2 + \lambda\epsilon \end{pmatrix} \\
&= \begin{pmatrix} 1.026 \\ -0.044 \end{pmatrix}
\end{aligned}$$

Furthemore, $\Delta\omega = \omega_c - \omega_b$ can be obtained.

$$\Delta\omega = \frac{1}{(1+\lambda)\epsilon^2 + \lambda(\lambda+2)} \begin{pmatrix} \epsilon^3 + \lambda\epsilon \\ -\epsilon^2 \end{pmatrix} = \begin{pmatrix} 0.0531 \\ -0.088 \end{pmatrix} \quad (12)$$

2.c

One can notice that $\Delta\omega$ with regularization is much smaller than $\Delta\omega$ without regularization. This implies that regularization makes the parameters less variable. This can be verified by the following figure.



3 LR with Regularization: A Probabilistic Perspective

4 Logistic Regression