

HW1- Theory

* Vectors are denoted by **boldfaced** characters, matrices by **BOLDFACED CAPITAL** letters.

1. Prove Normal Equations:

Given a training set $\mathcal{S} = \{\mathbf{X}, \mathbf{y}\}$, a linear hypothesis class $\{h_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{j=1}^N \theta_j x_j\}$ and the mean squared error loss function:

$$\mathcal{L} = \frac{1}{2M} \sum_{i=1}^M (h_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2$$

prove that $\boldsymbol{\theta}$ that minimizes \mathcal{L} satisfies:

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^T \mathbf{y}$$

where: $\mathbf{x}_i, \boldsymbol{\theta} \in \mathbb{R}^N$, $\mathbf{y} \in \mathbb{R}^M$, $\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ & \vdots & \\ - & \mathbf{x}_M^T & - \end{bmatrix}$, $M \geq N$

2. Unique solution:

Show that a unique solution for linear regression exists iff the features are not linearly dependent. Namely, show that a unique solution:

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathcal{L} = \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

exists iff \mathbf{X} has full column rank.

solution:

Guy Ilan

$$1. S = \{x, y\}$$

$$\text{hypothesis: } h_{\theta}(x) = \sum_{j=1}^N \theta_j x_j = x^T \theta$$

$$\text{loss function: } L = \frac{1}{2M} \sum_{i=1}^M (h_{\theta}(x_i) - y_i)^2$$

Proving that minimizing the MSE loss leads to $X^T X \theta = X^T y$:

Replace $h_{\theta}(x) = x^T \theta$:

$$L = \frac{1}{2M} \sum_{i=1}^M (x_i^T \theta - y_i)^2 \rightarrow L = \frac{1}{2M} (X\theta - y)^T (X\theta - y)$$

\rightarrow Minimizing the loss: $\frac{\partial L}{\partial \theta} = 0$

$$L = \frac{1}{2M} (X\theta - y)^T (X\theta - y) \rightarrow \frac{1}{2M} (\theta^T X^T - y^T) (X\theta - y) \rightarrow$$

$$\rightarrow \frac{1}{2M} (\theta^T X^T X \theta - \theta^T X^T y - y^T X \theta + y^T y)$$

$$\frac{\partial}{\partial \theta} (\theta^T X^T X \theta) = 2X^T X \theta, \quad \frac{\partial}{\partial \theta} (-\theta^T X^T y) = -X^T y, \quad \frac{\partial}{\partial \theta} (y^T y) = 0$$

$$\rightarrow \frac{\partial L}{\partial \theta} = \frac{1}{2M} (2X^T X \theta - 2X^T y) \rightarrow \frac{\partial L}{\partial \theta} = \frac{1}{M} (X^T X \theta - X^T y)$$



finding min by setting gradient equal to zero:

$$\frac{1}{N} (X^T X \theta - X^T y) = 0 \quad \cdot N \rightarrow (X^T X \theta - X^T y) = 0$$

$$\rightarrow X^T X \theta = X^T y$$

thus minimizing θ satisfies: $X^T X \theta = X^T y$

2. Show unique solution for normal equation:

$$X^T X \theta = X^T y \rightarrow \text{assuming that inverse of } X^T X \text{ exists: } X^T X \theta = X^T y \quad \cdot (X^T X)^{-1}$$

$$\rightarrow \theta = (X^T X)^{-1} X^T y$$

for $(X^T X)^{-1}$ to exist it must be invertible:

$X^T X$ is a $N \times N$ matrix and is invertible if and only if it has full rank.

By rank property $\text{rank}(X^T X) = \text{rank}(X)$.

Thus $X^T X$ is invertible if X has full rank.

X will have a full rank if all of its features are linearly independent.

$\text{rank}(X) = N$, implying $\text{rank}(X^T X) = N$

If X do not have a full rank, then $X^T X$ is not invertible and the normal equation will not have a unique solution:

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