

Deep Learning HW1

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

Solution:

Let us derive \mathcal{L} w.r to θ .

$$\mathcal{L} = \frac{1}{2M} \sum_{i=1}^M (h_{\theta}(x_i) - y_i)^2 = \frac{1}{2M} \sum_{i=1}^M \left(\sum_j^N \theta_j x_{ij} - y_i \right)^2 = \frac{1}{2M} \sum_{i=1}^M \|\theta X_i^T - y_i\|^2 = \frac{1}{2} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = 0 \text{ iff } \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) = 0$$

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

As required.

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:

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

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

Assume X is full rank, we know that $rk(X) = rk(X^T)$ and since $X^T X$ is square (matrix multiplication definition) That would suggest that $X^T X$ is also invertible since it's a composition of full rank transformation.

Taking our optimality condition from Q1 that states that $X^T X \theta = X^T y$ and applying our new found knowledge we could rearrange the terms to get $(X^T X)^{-1} X^T X \theta = (X^T X)^{-1} X^T y$ iff $(X^T X)^{-1} X^T y = \theta$. Meaning that the optimal point of \mathcal{L} with respect to θ is when $\theta = (X^T X)^{-1} X^T y$.

Assume that a unique solution $\theta = (X^T X)^{-1} X^T y$ exists. We know from basic linear algebra that $\operatorname{rank} X = \operatorname{rank} X^T X$ and since $X^T X$ is invertible, it is full rank, and so is X .

Short proof that $rk(X) = rk(X^T X)$. Proving that $\operatorname{Null}(X^T X) \subset \operatorname{Null}(X)$ suffices since the other direction is trivial, and the statement will follow.

Let u be a vector, such that $X^T X u = 0$ and let $v = Xu$. $u^T \overbrace{X^T X u}^0 = v^T v = 0$. That could only be if v is the 0 vector. And so $\operatorname{Null}(X^T X) \subset \operatorname{Null}(X)$.