

MATH REVIEW (FALL, 2020)

PROJECT

1. (10 points) Consider Linear Regression Model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where we observe data $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^p$. You may view \mathbf{X} and \mathbf{y} as a constant matrix and vector respectively. We assume $p \leq n$ and \mathbf{X} has full rank, that is, $\text{rank}(\mathbf{X}) = \min\{n, p\} = p$.

- (a) (3 points) The ordinary least squares estimator (LSE) $\boldsymbol{\beta}^{(LSE)}$ is the minimizer to ℓ_2 error

$$L_1(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

Find $\boldsymbol{\beta}^{(LSE)}$ and show it is the global minimizer using second derivative test.

- (b) (3 points) The ridge regression estimator $\boldsymbol{\beta}^{(ridge)}$ is the minimizer to ℓ_2 error with ℓ_2 regularization

$$L_2(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2,$$

where $\lambda > 0$ is a constant. Find the ridge regression estimator $\boldsymbol{\beta}^{(ridge)}$ and show it is the global minimizer using second derivative test.

- (c) (4 points) $\mathbf{X}\boldsymbol{\beta}^{(OLS)}$ is viewed as a projection of \mathbf{y} on the column space of \mathbf{X} . Let the orthogonal projection matrix be $P_{\mathbf{X}}$, write down its formula using \mathbf{X} and verify that it is a orthogonal projection matrix, that is, for any $\mathbf{u} \in \mathbb{R}^n$, $\langle P_{\mathbf{X}}\mathbf{u}, \mathbf{u} - P_{\mathbf{X}}\mathbf{u} \rangle = 0$.

2. (10 points) Consider a simple one-layer feedforward neural network model

$$\mathbf{y}_i = f_{NN}(\mathbf{x}_i) + \boldsymbol{\varepsilon}_i,$$

where $i = 1, \dots, n$, $\mathbf{x}_i \in \mathbb{R}^p$, $\mathbf{y}_i \in \mathbb{R}^d$ and

$$f_{NN}(\mathbf{x}_i) = B\phi(W\phi(A\mathbf{x}_i)).$$

The matrices $A \in \mathbb{R}^{m \times p}$, $B \in \mathbb{R}^{d \times m}$ and $W \in \mathbb{R}^{m \times m}$. The activation function ϕ is ReLU activation, that is, $\phi(\mathbf{u}) = (u_1 \vee 0, u_2 \vee 0, \dots, u_n \vee 0)^\top$ for $\mathbf{u} \in \mathbb{R}^n$. Thus, you may write $D_{i,0}A\mathbf{x}_i = \phi(A\mathbf{x}_i)$ and $D_{i,1}W\mathbf{h}_{i,0} = \phi(W\mathbf{h}_{i,0})$ where $D_{i,0}$ and $D_{i,1}$ are diagonal matrices with elements being 0 or 1. Consider the ℓ_2 error

$$L(W) = \frac{1}{2} \sum_{i=1}^n \|\mathbf{y}_i - f_{NN}(\mathbf{x}_i)\|_2^2.$$

Find the gradient matrix of $L(W)$, $\nabla_W L = \left(\frac{\partial L(W)}{\partial W_{s,t}} \right)_{m \times m}$. Be careful with dimensions when you work on vectors and matrices.

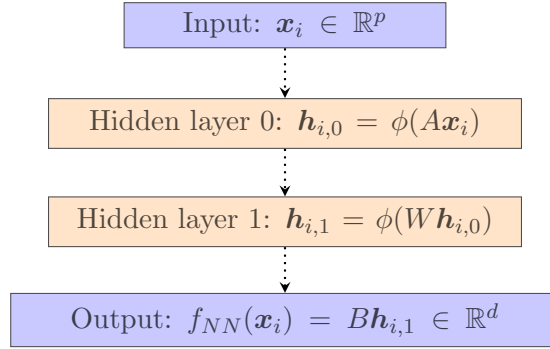


Figure 1: One-layer Neural Network