

HW2

Thursday, February 10, 2022 11:25 PM

1. [10 points] In Module 2, we gave the normal equation (i.e., closed-form solution) for linear regression using MSE as the cost function. **Prove that the closed-form solution for Ridge Regression**

is $\mathbf{w} = (\lambda I + \mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$, where I is the identity matrix, $\mathbf{X} = (x^{(1)}, x^{(2)}, \dots, x^{(m)})^T$ is the input data matrix, $x^{(i)} = (1, x_1, x_2, \dots, x_n)$ is the i -th data sample, and $\mathbf{y} = (y^{(1)}, y^{(2)}, \dots, y^m)$. Assume the hypothesis function $h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$, and $y^{(j)}$ is the measurement of $h_{\mathbf{w}}(\mathbf{x})$ for the j -th training sample. The cost function of the Ridge Regression is $E(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^T \cdot \mathbf{x}^{(i)} - y^{(i)})^2 + \lambda \sum_{i=1}^m w_i^2$.

Ridge Regression Closed Form: $\mathbf{w} = (\lambda I + \mathbf{X}^T \cdot \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$

Ridge Regression Cost Function: $E(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^T \cdot \mathbf{x}^{(i)} - y^{(i)})^2 + \lambda \sum_{i=1}^m w_i^2$

Shape of $\mathbf{X} \rightarrow m \times n+1$, Shape of $\mathbf{y} \rightarrow m$ (column vector)

Shape of $\mathbf{w} \rightarrow n+1$ (column vector), Shape of $\mathbf{x} \rightarrow n+1$ (column vector)

$$\begin{aligned} E(\mathbf{w}) &= \sum_{i=1}^m (\mathbf{w}^T \cdot \mathbf{x}^{(i)} - y^{(i)})^T \cdot (\mathbf{w}^T \cdot \mathbf{x}^{(i)} - y^{(i)}) + \lambda \sum_{i=1}^m w_i^T \cdot w_i \\ &= \sum_{i=1}^m (\mathbf{x}^{(i)T} \cdot \mathbf{w} - y^{(i)T}) \cdot (\mathbf{w}^T \cdot \mathbf{x}^{(i)} - y^{(i)}) + \lambda \sum_{i=1}^m w_i^T \cdot w_i \\ &= \sum_{i=1}^m \mathbf{x}^{(i)T} \cdot \mathbf{w} \cdot \mathbf{w}^T \cdot \mathbf{x}^{(i)} - \mathbf{x}^{(i)T} \cdot \mathbf{w} \cdot \mathbf{y}^{(i)} - \mathbf{y}^{(i)T} \cdot \mathbf{w}^T \cdot \mathbf{x}^{(i)} + \mathbf{y}^{(i)T} \cdot \mathbf{y}^{(i)} + \lambda \sum_{i=1}^m w_i^T \cdot w_i \\ &= \sum_{i=1}^m \mathbf{x}^{(i)T} \cdot \mathbf{w} \cdot \mathbf{w}^T \cdot \mathbf{x}^{(i)} - 2 \mathbf{w}^T \cdot \mathbf{x}^{(i)} \cdot \mathbf{y}^{(i)} + \mathbf{y}^{(i)T} \cdot \mathbf{y}^{(i)} + \lambda \sum_{i=1}^m w_i^T \cdot w_i \\ &= \mathbf{w}^T [\sum_{i=1}^m \mathbf{x}^{(i)T} \cdot \mathbf{x}^{(i)}] \cdot \mathbf{w} - 2 \mathbf{w}^T \sum_{i=1}^m \mathbf{x}^{(i)} \cdot \mathbf{y}^{(i)} + \sum_{i=1}^m \mathbf{y}^{(i)T} \cdot \mathbf{y}^{(i)} + \lambda \sum_{i=1}^m w_i^T \cdot w_i \\ &= \mathbf{w}^T \cdot \mathbf{X} \cdot \mathbf{X}^T \cdot \mathbf{w} - 2 \mathbf{w}^T \mathbf{X} \cdot \mathbf{y} + \mathbf{y}^T \cdot \mathbf{y} + \lambda \mathbf{w}^T \cdot \mathbf{w} \end{aligned}$$

$$\nabla E(\mathbf{w}) = 2 \mathbf{X} \cdot \mathbf{X}^T \mathbf{w} - 2 \mathbf{X} \cdot \mathbf{y} + 2 \lambda \mathbf{w} = 0 \Rightarrow \text{set to 0 to optimize}$$

$$\rightarrow 2 \mathbf{X} \cdot \mathbf{X}^T \mathbf{w} + 2 \lambda \mathbf{w} = 2 \mathbf{X}^T \mathbf{y}$$

$$\rightarrow \mathbf{X} \cdot \mathbf{X}^T \mathbf{w} + \lambda \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\rightarrow \mathbf{w} (\mathbf{X} \cdot \mathbf{X}^T + \lambda I) = \mathbf{X}^T \mathbf{y} \Rightarrow \boxed{\mathbf{w} = (\mathbf{X} \cdot \mathbf{X}^T + \lambda I)^{-1} \cdot \mathbf{X}^T \mathbf{y}}$$

2. [10 points] Assume we have K different classes in a multi-class Softmax Regression model. The

posterior probability is $\hat{p}_k = \delta(s_k(x))_k = \frac{\exp(s_k(x))}{\sum_{j=1}^K \exp(s_j(x))}$ for $k = 1, 2, \dots, K$, where $s_k(x) = \theta_k^T \cdot x$,

input x is an n -dimension vector, and K the total number of classes.

- 1) To learn this Softmax Regression model, how many parameters we need to estimate? What are these parameters?

- 2) Consider the cross-entropy cost function $J(\theta)$ of m training samples $\{(x_i, y_i)\}_{i=1,2,\dots,m}$ as below. Derive the gradient of $J(\theta)$ regarding to θ_k .

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)})$$

where $y_k^{(i)} = 1$ if the i th instance belongs to class k ; 0 otherwise.

softmax function

Posterior Probability $\hat{p}_k = \delta(s_k(x))_k = \frac{e^{s_k(x)}}{\sum_{j=1}^k e^{s_j(x)}}$, $s_k(x) = \theta_k^T \cdot x$
 Shape of $x \rightarrow n$ (column vector), $K = \#$ of classes

1) We have that $s_k(x) = \theta_k^T \cdot x$ & x is a n -dimension vector.
 This means that each θ_k is an n -dimension vector. Since there are K classes, we need to estimate $n \times K$ parameters.
 The parameters θ_k are vectors of parameters that are used to transform the input feature vectors x to $s_k(x)$, after which the softmax is applied to get a final prediction.

2) Cost Function : $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)})$

$$\nabla J(\theta) = -\frac{1}{m} \sum_{i=1}^m \nabla \sum_{k=1}^K y_k \cdot \log(\hat{p}_k)$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K y_k \cdot \nabla \log(\hat{p}_k) \right] = -\frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K \frac{y_k}{\hat{p}_k} \cdot \frac{\partial \hat{p}_k}{\partial s_k} \right]$$

Need to calculate derivative of softmax : $\frac{\partial}{\partial s_i} \log(\hat{p}_k)$ log derivative to simplify

$$\frac{\partial}{\partial s_i} \log(\hat{p}_k) = \frac{1}{\hat{p}_k} \cdot \frac{\partial \hat{p}_k}{\partial s_i} * \frac{\partial}{\partial x} \log(x) = \frac{1}{x} * \text{chain rule}$$

$$\frac{\partial \hat{p}_k}{\partial s_i} = \hat{p}_k \cdot \frac{\partial}{\partial s_i} \log(\hat{p}_k) * \text{rearrange}$$

$$= \hat{p}_k \cdot \frac{\partial}{\partial s_i} \log\left(\frac{e^{s_k}}{\sum_{k=1}^K e^{s_k}}\right) = \hat{p}_k \cdot \left[\frac{\partial}{\partial s_i} \left(s_k - \log\left(\sum_{k=1}^K e^{s_k}\right) \right) \right] * \text{substitute } \hat{p}_k$$

$$= \hat{p}_k \cdot \frac{\partial s_k}{\partial s_i} - \frac{\partial}{\partial s_i} \log\left(\sum_{k=1}^K e^{s_k}\right) * \text{distribute}$$

$$= \hat{p}_k \cdot \left[\left\{ \{k=j\} - \frac{\partial}{\partial s_i} \log\left(\sum_{k=1}^K e^{s_k}\right) \right\} \right] * \frac{\partial s_k}{\partial s_i} = \begin{cases} 1, & \text{if } k=j \\ 0, & \text{otherwise} \end{cases}$$

$$= \hat{p}_k \cdot \left[\left\{ \{k=j\} - \frac{1}{\sum_{k=1}^K e^{s_k}} \cdot \left(\frac{\partial}{\partial s_i} \sum_{k=1}^K e^{s_k} \right) \right\} \right] * \text{chain rule}$$

$$= \hat{p}_k \cdot \left[\left\{ \{k=j\} - \frac{e^{s_i}}{\sum_{k=1}^K e^{s_k}} \right\} \right] * \text{partial derivative of } \sum = e^{s_i}$$

$$= \hat{p}_k \cdot \left[\left\{ \{k=j\} - p_i \right\} \right] * \text{substitute } p_i \text{ for posterior}$$

Now plug derivative of softmax back into $\nabla J(\theta)$

$$\nabla J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K y_k \cdot \left(\left\{ \{k=j\} - p_i \right\} \right) \right] \cdot X^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K y_k \cdot \left[\left\{ \{k=j\} - \sum_{k=1}^K y_k \cdot p_i \right\} \right] \right] \cdot X^{(i)} * \text{distribute}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y_k - \sum_{k=1}^K y_k \cdot p_i \right] \cdot X^{(i)} * \sum \left\{ \{k=j\} \right\} = 1$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y_k - p_i \sum_{k=1}^K y_k \right] \cdot X^{(i)} * p_i \text{ comes out of } \sum$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y_k - p_i \cdot (1) \right] \cdot X^{(i)} * \sum y_k = 1 \text{ (one-hot)}$$

$$\Rightarrow \boxed{-\frac{1}{m} \sum_{i=1}^m \left[\hat{p}_k^{(i)} - y_k^{(i)} \right] \cdot X^{(i)}}$$