CSC 4792 - DATA MINING

PROGRAMMING EXERCISE 1//2019

Question One

Given the data matrix $\mathbf{X} \in \mathbb{R}^{Mx(d+1)}$ whose rows are the inputs $x_i$ as row vectors, and target vector $y \in \mathbb{R}^{Mx1}$ a column vector whose components are the target values $y_i$. In class we discussed that in order to derive the linear regression algorithm, we need to minimize the in-sample error function of $w$ and the data $\mathbf{X}, \mathbf{y}$ given by $E_{in}(w) = \frac{1}{M}\sum_{i=1}^{M}(w^T x_i - y_i)^2$, over all possible $w \in \mathbb{R}^{d+1}$, the optimization problem formalized as $w = argmin_{w \in \mathbb{R}^{d+1}} E_{in}(w)$. Show that besides gradient descent method, there is a closed-form solution to the problem given by $\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$, assuming $\mathbf{X}^T\mathbf{X}$ is invertible.

Question Two

Consider the following data sampled from employees of a given organization about age against their income.

| ID | age | income |
|----|-----|--------|
| 1 | 21 | 24000 |
| 2 | 32 | 48000 |
| 3 | 62 | 83000 |
| 4 | 72 | 61000 |
| 5 | 84 | 52000 |

Tasks:

i.      Plot the datapoints (age against income) to visualize the distribution in 2D.

ii.     Using the closed-form solution in *question 1* or gradient descent, find the model that *suitably* predicts the income given the age of the employee in this organization.

iii.    Plot your model in Question 2(ii) to visualize how it fits the datapoints.

iv.     Using your model in Question (ii), what is the expected income if the employee is 50 year old?

v.      If your model does not fit the dataset properly, use Polynomial regression to find the best fit model.

Prepared By: Claytone Sikasote