

Solution for Assignment 1

Chinh La

2018-09-25

1 Newton's method for computing least squares

- (a) Find the Hessian of cost function $J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$

Answer:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 \quad (1)$$

$$\frac{\partial J(\theta)}{\partial \theta_i} = \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x^{(i)} \quad (2)$$

$$\frac{\partial^2 J(\theta)}{\partial \theta_i \partial \theta_j} = \sum_{i=1}^m x^{(j)} x^{(i)} \quad (3)$$

The **Hessian** matrix of cost function $J(\theta)$ is an $n \times n$ matrix where

$$H_{ij} = \sum_{i=1}^m x^i x^j = X^T X_{ji}$$

with $X = [x^{(1)} \quad x^{(2)} \quad \dots \quad x^{(m)}]^T$

- (b) Show that the first iteration of Newton's method gives us $\theta^* = (X^T X)^{-1} X^T \vec{y}$

Answer: As the notes 1 show, $\theta \in R^{1 \times n}$, Newton's method performs the following update:

$$\theta_{i+1} := \theta_i - H^{-1} \nabla_{\theta} l(\theta_i) \quad (4)$$

In this problem, $\nabla_{\theta} l(\theta) = X^T (X\theta^T - \vec{y})$, $\vec{y} = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]^T$, $H^{-1} = (X^T X)^{-1}$.

Let $\theta_0 = \vec{0}^T$, we have

$$\Rightarrow \theta_1 = \theta_0 - (X^T X)^{-1} X^T (X\theta^T - \vec{y}) \quad (5)$$

$$\Rightarrow \theta_1 = \vec{0}^T - (X^T X)^{-1} X^T (X\vec{0} - \vec{y}) \quad (6)$$

$$\Rightarrow \theta_1 = (X^T X)^{-1} X^T \vec{y} \quad (7)$$

So θ_1 is optimal solution of cost function $l(\theta)$ and we only need 1 iteration.

2 Locally-weighted logistic regression

See alq2.py file

3 Multivariate least squares

(a) The cost function for this case is

$$J(\Theta) = \frac{1}{2} \sum_{n=1}^m \sum_{j=1}^p \left(\left(\Theta^T x^{(i)} \right)_j - y_j^{(i)} \right)^2$$

Find the matrix way to write it.

Answer:

$$\text{Let } X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix}, Y = \begin{bmatrix} y_1^{(1)} & y_2^{(2)} & \cdots & y_p^{(1)} \\ y_1^{(2)} & y_2^{(2)} & \cdots & y_p^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ y_p^{(m)} & y_2^{(m)} & \cdots & y_p^{(m)} \end{bmatrix}, \Theta = \begin{bmatrix} \theta_{1(1)} & \theta_{1(2)} & \cdots & \theta_{1(p)} \\ \theta_{2(1)} & \theta_{2(2)} & \cdots & \theta_{2(p)} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{n(1)} & \theta_{n(2)} & \cdots & \theta_{n(p)} \end{bmatrix}$$

We have

$$J(\Theta) = \text{tr}((X\Theta - Y)^T(X\Theta - Y)) \quad (8)$$

$$J(\Theta) = \sum_{i=1}^m ((X\Theta - Y)^T(X\Theta - Y))_{ii} \quad (9)$$

$$J(\Theta) = \sum_{i=1}^m ((X\Theta - Y)^T(X\Theta - Y))_{ii} \quad (10)$$

where $((X\Theta - Y)^T(X\Theta - Y))_{ii}$ is the i th element on main diagonal of $(\Theta^T X - Y^T)(\Theta^T X - Y)$ and

$$(x\Theta - Y) = [\theta^T x^{(1)} - \theta^T x^{(2)} \quad \cdots \quad \theta^T x^{(m)}] \quad (11)$$

$$((X\Theta - Y)^T(X\Theta - Y))_{ii} = \sum_{j=1}^p \left(\left(\Theta^T x^{(i)} \right)_j - y_j^{(i)} \right)^2 \quad (12)$$

$$(10)(12) \Rightarrow J(\Theta) = \frac{1}{2} \sum_{n=1}^m \sum_{j=1}^p \left(\left(\Theta^T x^{(i)} \right)_j - y_j^{(i)} \right)^2 \quad (13)$$

(b) Find the closed form solution for Θ which minimizes $J(\Theta)$. This is equivalent to the normal equations for multivariate case.

Answer:

$$\nabla_{\Theta} J(\Theta) = \frac{1}{2} \nabla_{\Theta} \text{tr}((X\Theta - Y)^T(X\Theta - Y)) \quad (14)$$

$$\Rightarrow \nabla_{\Theta} J(\Theta) = \frac{1}{2} \nabla_{\Theta} \text{tr}((\Theta^T X^T - Y^T)(X\Theta - Y)) \quad (15)$$

$$\Rightarrow \nabla_{\Theta} J(\Theta) = \frac{1}{2} \nabla_{\Theta} \text{tr}(\Theta^T X^T X\Theta - \Theta^T X^T Y - Y^T X\Theta + Y^T Y) \quad (16)$$

$$\Rightarrow \nabla_{\Theta} J(\Theta) = \frac{1}{2} \nabla_{\Theta} \text{tr}(\Theta^T X^T X\Theta) - \nabla_{\Theta} \text{tr}(\Theta^T X^T Y) - \nabla_{\Theta} \text{tr}(Y^T X\Theta) \quad (17)$$

$$\Rightarrow \nabla_{\Theta} J(\Theta) = \frac{1}{2} ((X^T X)^T \Theta + (X^T X)\Theta - X^T Y - (Y^T X)^T) \quad (18)$$

$$\Rightarrow \nabla_{\Theta} J(\Theta) = (X^T X)\Theta - X^T Y \quad (19)$$

Let $\nabla_{\Theta} J(\Theta) = 0$ and we have the optimal Θ^*

$$\Theta^* = (X^T X)^{-1} X^T Y \quad (20)$$

(c) Suppose instead of considering the multivariate vectors $y^{(i)}$ all at once, we instead compute each variable $y_j^{(i)}$ separately for each $j = 1, \dots, p$. In this case, we have a p individual linear models, of the form

$$y_j^{(i)} = \theta_j^T x^{(i)}, j = 1, \dots, p$$

(So here, each $\theta_j \in \mathbb{R}^n$). How do the parameters from these p independent least squares problems compare to the multivariate solution?

Answer:

As we have

$$\theta_i = (X^T X)^{-1} X^T \vec{y}_i \quad (21)$$

where $\Theta = [\theta_1 \ \theta_2 \ \cdots \ \theta_m]$, $Y = [y_1 \ y_2 \ \cdots \ y_m]$

Then the solution of each independent model become the solution of multivariate case.

4 Naive Bayes

In this problem, we look at maximum likelihood parameter estimation using the naive Bayes assumption. Here, the input features $x_j, j = 1, \dots, n$ to our model are discrete, binary-valued variable, so $x_j \in \{0, 1\}$. We call $x = [x_1 \ x_2 \ \dots \ x_n]$ to be the input vector. For each training example, our output targets are single binary-value $y \in \{0, 1\}$. Our model is then parameterized by $\phi_{j|y=0} = p(x_j = 1|y_j = 0)$, $\phi_{j|y=1} = p(x_j = 1|y_j = 1)$ and $\phi_y = p(y = 1)$. We model the joint distribution of (x, y) according to

$$\begin{aligned} p(y) &= (\phi_y)^y (1 - \phi_y)^{1-y} \\ p(x|y=0) &= \prod_{j=1}^n p(x_j|y=0) \\ &= \prod_{j=1}^n (\phi_{j|y=0})^{x_j} (1 - \phi_{j|y=0})^{1-x_j} \\ p(x|y=1) &= \prod_{j=1}^n p(x_j|y=1) \\ &= \prod_{j=1}^n (\phi_{j|y=1})^{x_j} (1 - \phi_{j|y=1})^{1-x_j} \end{aligned}$$

- (a) Find the joint likelihood function $l(\varphi) = \log \left(\prod_{i=1}^m p(x^{(i)}, y^{(i)}; \varphi) \right)$ in terms of the model parameters given above. Here, φ represents the entire parameters $\{\phi_y, \phi_{y=0}, \phi_{y=1}, j = 1, \dots, n\}$.

Answers:

$$l(\varphi) = \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \varphi) \quad (22)$$

$$l(\varphi) = \sum_{i=1}^m \log \left(p(x^{(i)}, y^{(i)}; \varphi) \right) \quad (23)$$

$$l(\varphi) = \sum_{i=1}^m \log \left(p(x^{(i)}|y^{(i)}; \varphi) p(y^{(i)}; \varphi) \right) \quad (24)$$

$$l(\varphi) = \sum_{i=1}^m \left(\log p(x^{(i)}|y^{(i)}; \varphi) + \log p(y^{(i)}; \varphi) \right) \quad (25)$$

$$l(\varphi) = \sum_{i=1}^m \left(\log \prod_{j=1}^n p(x_j^{(i)}|y^{(i)}; \varphi) + \log p(y^{(i)}; \varphi) \right) \quad (26)$$

$$l(\varphi) = \sum_{i=1}^m \left(\sum_{j=1}^n \log p(x_j^{(i)}|y^{(i)}; \varphi) + \log p(y^{(i)}; \varphi) \right) \quad (27)$$

$$l(\varphi) = \sum_{i=1}^m \left(\sum_{j=1}^n \left(x_j^{(i)} \log \phi_{j|y} + (1 - x_j^{(i)}) \log(1 - \phi_{j|y}) \right) + y^{(i)} \log \phi_y + (1 - y^{(i)}) \log(1 - \phi_y) \right) \quad (28)$$

- (b) Show that the parameters which maximize the likelihood function are the same as those given in the lecture notes: i.e., that

$$\begin{aligned}\phi_{j|y=0} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \\ \phi_{j|y=1} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \\ \phi_y &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m}\end{aligned}$$

Answers:

$$\nabla_{\phi_{j|y=0}} l(\varphi) = \sum_{i=1}^m \left(\frac{x_j^{(i)}}{\phi_{j|y=0}} 1\{y^{(i)} = 0\} - \frac{1 - x_j^{(i)}}{1 - \phi_{j|y=0}} 1\{y^{(i)} = 0\} \right) \quad (29)$$

$$\Rightarrow \nabla_{\phi_{j|y=0}} l(\varphi) = 0, \sum_{i=1}^m \left(\frac{x_j^{(i)}}{\phi_{j|y=0}} - \frac{1 - x_j^{(i)}}{1 - \phi_{j|y=0}} \right) 1\{y^{(i)} = 0\} = 0 \quad (30)$$

$$\Rightarrow \sum_{i=1}^m \left(x_j^{(i)} (1 - \phi_{j|y=0}) - (1 - x_j^{(i)}) \phi_{j|y=0} \right) 1\{y^{(i)} = 0\} = 0 \quad (31)$$

$$\Rightarrow \sum_{i=1}^m \left(x_j^{(i)} - \phi_{j|y=0} \right) 1\{y^{(i)} = 0\} = 0 \quad (32)$$

$$\Rightarrow \phi_{j|y=0} = \frac{\sum_{i=1}^m x_j^{(i)} 1\{y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \quad (33)$$

$$\Rightarrow \phi_{j|y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \quad (34)$$

With $\nabla_{\phi_{j|y=1}} l(\varphi) = 0$ we have the same result

$$\phi_{j|y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \quad (35)$$

With $\nabla_{\phi_y} l(\varphi) = 0$ we have:

$$\nabla_{\phi_y} l(\varphi) = \sum_{i=1}^m \left(\frac{y^{(i)}}{\phi_y} - \frac{1 - y^{(i)}}{1 - \phi_y} \right) \quad (36)$$

$$\sum_{i=1}^m \left(y^{(i)} (1 - \phi_y) - (1 - y^{(i)}) \phi_y \right) = 0 \quad (37)$$

$$\sum_{i=1}^m \left(y^{(i)} - \phi_y \right) = 0 \quad (38)$$

$$\sum_{i=1}^m y^{(i)} - m\phi_y = 0 \quad (39)$$

$$\phi_y = \frac{\sum_{i=1}^m y^{(i)}}{m} = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m} \quad (40)$$

- (c) Consider making a prediction on some new data point x using the most likely class estimate generated by the naive Bayes algorithm. Show that the hypothesis returned by naive Bayes is a linear classifier - i.e., if $p(y = 0|x)$ and $p(y = 1|x)$ are the class probabilities returned by naive Bayes, show that there exists some $\theta \in \mathbb{R}^{n+1}$ such that

$$p(y = 1|x) \geq p(y = 0|x) \text{ if and only if } \theta^T \begin{bmatrix} 1 \\ x \end{bmatrix} \geq 0$$

(Assume θ_0 is an intercept terms.)

Answers:

$$p(y = 1|x) \geq p(y = 0|x) \quad (41)$$

$$\iff \frac{p(y = 1|x)}{p(y = 0|x)} \geq 1 \quad (42)$$

$$\iff \frac{\left(\prod_{j=1}^n p(x_j|y = 1)\right) p(y = 1)}{\left(\prod_{j=1}^n p(x_j|y = 0)\right) p(y = 0)} \geq 1 \quad (43)$$

$$\iff \frac{\left(\prod_{j=1}^n (\phi_{j|y=1})^{x_j} (1 - \phi_{j|y=1})^{1-x_j}\right) \phi_y}{\left(\prod_{j=1}^n (\phi_{j|y=0})^{x_j} (1 - \phi_{j|y=0})^{1-x_j}\right) (1 - \phi_y)} \geq 1 \quad (44)$$

$$\iff \sum_{j=1}^n \left(x_j \log \left(\frac{\phi_{j|y=1}}{\phi_{j|y=0}} \right) + (1 - x_j) \log \left(\frac{1 - \phi_{j|y=1}}{1 - \phi_{j|y=0}} \right) \right) + \log \left(\frac{\phi_y}{1 - \phi_y} \right) \geq 0 \quad (45)$$

$$\iff \sum_{j=1}^n x_j \log \left(\frac{(\phi_{j|y=1})(1 - \phi_{j|y=0})}{(\phi_{j|y=0})(1 - \phi_{j|y=1})} \right) + \sum_{j=1}^n \log \left(\frac{1 - \phi_{j|y=1}}{1 - \phi_{j|y=0}} \right) + \log \left(\frac{\phi_y}{1 - \phi_y} \right) \geq 0 \quad (46)$$

$$\iff \theta^T \begin{bmatrix} 1 \\ x \end{bmatrix} \geq 0, \quad (47)$$

where

$$\begin{aligned} \theta_0 &= \sum_{j=1}^n \log \left(\frac{1 - \phi_{j|y=1}}{1 - \phi_{j|y=0}} \right) + \log \left(\frac{\phi_y}{1 - \phi_y} \right) \\ \theta_j &= \log \left(\frac{(\phi_{j|y=1})(1 - \phi_{j|y=0})}{(\phi_{j|y=0})(1 - \phi_{j|y=1})} \right), j = 1, \dots, n \end{aligned}$$