
Machine Learning, 2024 Spring

Assignment 3

Notice

Plagiarizer will get 0 points.
L^AT_EX is highly recommended. Otherwise you should write as legibly as possible.

Problem 1 For logistic regression, show that

$$\nabla E_{\text{in}}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}} = \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n) \quad (1)$$

Argue that a ‘misclassified’ example contributes more to the gradient than a correctly classified one.

Solution

1. Let $\theta(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$ be the sigmoid function.

So we have $p(y_n | \mathbf{x}_n; \mathbf{w}) = \theta(y_n \mathbf{w}^T \mathbf{x}_n)$, and the maximum likelihood estimation is

$$\max \prod_{n=1}^N p(y_n | \mathbf{x}_n; \mathbf{w})$$

For a convenient calculation, we take the log of the likelihood function, and then take the negative of it, and take mean of each term, so we have what we want to minimize:

$$\begin{aligned} E_{\text{in}}(\mathbf{w}) &= -\frac{1}{N} \log \left[\prod_{n=1}^N p(y_n | \mathbf{x}_n; \mathbf{w}) \right] \\ &= -\frac{1}{N} \sum_{n=1}^N \log \theta(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= -\frac{1}{N} \sum_{n=1}^N \log \frac{1}{1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}} \\ &= -\frac{1}{N} \sum_{n=1}^N [-\log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})] \\ &= \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) \end{aligned}$$

So the gradient of $E_{\text{in}}(\mathbf{w})$ is

$$\begin{aligned}
\nabla E_{\text{in}}(\mathbf{w}) &= \nabla \left[\frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) \right] \\
&= \frac{1}{N} \sum_{n=1}^N \nabla \log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) \\
&= \frac{1}{N} \sum_{n=1}^N \frac{(-y_n \mathbf{x}_n) e^{-y_n \mathbf{w}^T \mathbf{x}_n}}{1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}} \\
&= -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}} \\
&= \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n)
\end{aligned}$$

2. From the property of the sigmoid function, we know that the sigmoid function is monotonic increasing.

So for a correctly classified example, we have $y_n \mathbf{w}^T \mathbf{x}_n \geq 0$, so $\theta(-y_n \mathbf{w}^T \mathbf{x}_n) \in (0, 0.5]$.

And for a misclassified example, we have $y_n \mathbf{w}^T \mathbf{x}_n < 0$, so we have $\theta(-y_n \mathbf{w}^T \mathbf{x}_n) \in (0.5, 1)$.

From the gradient formula $\nabla E_{\text{in}}(\mathbf{w})$, the difference between the correctly classified example and the misclassified example is only the sigmoid function.

And we also have the misclassified example's $\theta(-y_n \mathbf{w}^T \mathbf{x}_n)$ is larger than the correctly classified example's $\theta(-y_n \mathbf{w}^T \mathbf{x}_n)$, so the misclassified example contributes more to the gradient than a correctly classified one.

So above all, we have proved that

$$\nabla E_{\text{in}}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}} = \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n)$$

And we also proved that the 'misclassified' example contributes more to the gradient than a correctly classified one.

Problem 2 For linear regression, the out-of-sample error is

$$E_{\text{out}}(h) = \mathbb{E}[(h(x) - y)^2]. \quad (2)$$

Show that among all hypotheses, the one that minimizes E_{out} is given by

$$h^*(x) = \mathbb{E}[y|x]. \quad (3)$$

The function h^* can be treated as a deterministic target function, in which case we can write $y = h^*(x) + \epsilon(x)$ where $\epsilon(x)$ is an (input dependent) noise variable. Show that $\epsilon(x)$ has expected value zero.

Solution

Suppose the joint distribution of x and y is $f(x, y)$, the distribution of x is $f_X(x)$, the distribution of y is $f_Y(y)$, and the conditional distribution of y given x is $f_{Y|X}(y|x)$.

Lemma: The law of iterated Expectation

$$\mathbb{E}(Y) = \mathbb{E}_X[\mathbb{E}(Y|X)]$$

proof:

$$\begin{aligned} \mathbb{E}_X[\mathbb{E}(Y|X)] &= \mathbb{E}_X \left[\int_{-\infty}^{+\infty} y f_{Y|X}(y|x) dy \right] \\ &= \mathbb{E}_X \left[\int_{-\infty}^{+\infty} y \frac{f(x, y)}{f_X(x)} dy \right] \\ &= \int_{-\infty}^{+\infty} \left[\int_{-\infty}^{+\infty} y \frac{f(x, y)}{f_X(x)} dy \right] f_X(x) dx \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} y f(x, y) dx dy \\ &= \int_{-\infty}^{+\infty} y \left[\int_{-\infty}^{+\infty} f(x, y) dx \right] dy \\ &= \int_{-\infty}^{+\infty} y f_Y(y) dy \\ &= \mathbb{E}(Y) \end{aligned}$$

1. Since $E_{\text{out}}(h) = \mathbb{E}[(h(x) - y)^2] = \mathbb{E}_{X,Y}[(h(x) - y)^2]$, we have

$$\begin{aligned} \frac{\partial E_{\text{out}}(h)}{\partial h} &= 2\mathbb{E}_{X,Y}[h(x) - y] \\ &= 2\mathbb{E}_X[\mathbb{E}[(h(x) - y)|x]] \quad (\text{Lemma}) \\ &= 2\mathbb{E}_X[\mathbb{E}[h(x)|x] - \mathbb{E}[y|x]] \\ &= 2\mathbb{E}_X[h(x) - \mathbb{E}[y|x]] \end{aligned}$$

The necessary condition for minimizing the $E_{\text{out}}(h)$ is to let $2\mathbb{E}_X[h(x) - \mathbb{E}[y|x]] = 0$. Since $\mathbb{E}_X[h(x) - \mathbb{E}[y|x]]$ is the function of x , so if $h(x) - \mathbb{E}[y|x] = 0$, then $\frac{\partial E_{\text{out}}(h)}{\partial h} = 0$, which is the necessary condition that minimizes $E_{\text{out}}(h)$.

So we have $h^*(x) = \mathbb{E}[y|x]$.

2. Since $y = h^*(x) + \epsilon(x)$, we have

$$\mathbb{E}[\epsilon(x)] = \mathbb{E}[y - h^*(x)] = \mathbb{E}(y) - \mathbb{E}[h^*(x)] = \mathbb{E}[y] - \mathbb{E}[\mathbb{E}[y|x]] \stackrel{\text{Lemma}}{=} \mathbb{E}[y] - \mathbb{E}[y] = 0$$

So above all, we have prove that $h^*(x) = \mathbb{E}[y|x]$ and $\epsilon(x)$ has expected value zero.

Problem 3 According to the iterative format provided as follow:

- Use the SUV dataset to implement (using Python or MATLAB) the Gradient Descent method to find the optimal model for logistic regression.
- What is your test error? What are the sizes of the training set and test set, respectively?
- What is your learning rate? How was it chosen? How many steps were iterated in total?
- Present the results of the last 10 steps produced by your algorithm, including the loss, learning rate, the L2 norm of the gradient, and the number of function evaluations and gradient evaluations.

Fixed learning rate gradient descent:

- 1: Initialize the weights at time step $t = 0$ to $\mathbf{w}(0)$.
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: Compute the gradient $\mathbf{g}_t = \nabla E_{\text{in}}(\mathbf{w}(t))$.
- 4: Set the direction to move, $\mathbf{v}_t = -\mathbf{g}_t$.
- 5: Update the weights: $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$.
- 6: Iterate to the next step until it is time to stop.
- 7: Return the final weights.

Dataset reference and description

Dataset and download

Solution

0. Data preperation

After checking the SUV dataset, we could find that each buyer's information has a unique 'User ID', so we can just drop the 'User ID' column.

To use the 'Gender' information, we map 'Male' into 1, and 'Female' into 0.

If we use exponential function, it may be too big with large 'Age' and 'EstimatedSalary', so we can separate the 'Age' and 'EstimatedSalary' into different intervals. And use an one-hot vector to represent the which range the 'Age' and 'EstimatedSalary' belongs to. And the modified encoding for 'Age' and 'Estimated salary' are as followed in Figure 1 and Figure 2. The catagoies of 'Age' and 'EstimatedSalary' are chosen by the distribution graph of the original data, and each catagory is represented by a one-hot vector.

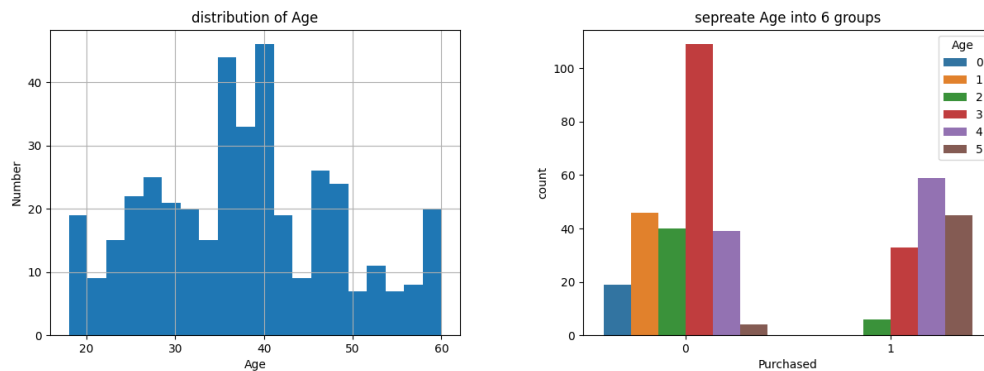


Figure 1: Analysis of Age

To suit better for logistic regression, the 'Purchased' column is also modified into '-1' for unpurchased and '1' for purchased.

And before sending the data into the logistic regression, the input features are normalized.

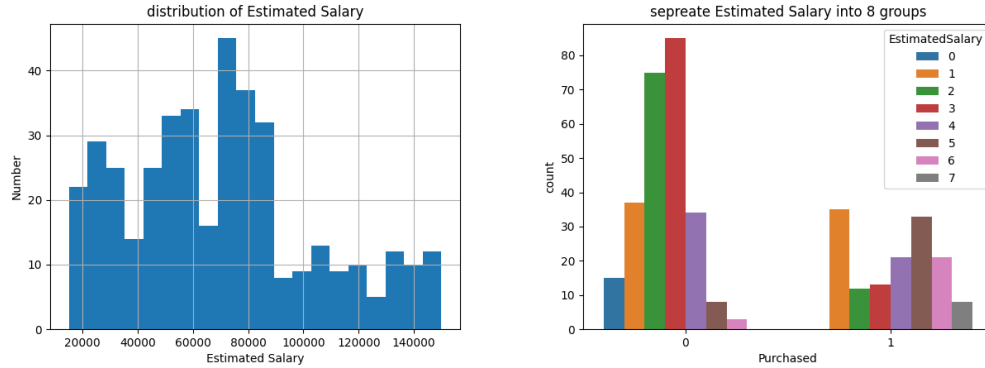


Figure 2: Analysis of Estimated Salary

Then the 'Gender', 'Age', 'EstimatedSalary' can be used as the input features, and the 'Purchased' can be used as the output label. So this is a binary classification problem. And we can use the logistic regression to solve this problem.

1. Since we are applying the logistic regression, so the testing error is:

$$E(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})$$

which have been proved in Problem 1. And N is the size of the testing set's size. The variation curve for the testing error is shown in Figure ????????????????

And the size of training set is 80 percent, which is totally 280 samples. And the testing set is 30 percent, which is totally 120 samples. And before separating the data, we random shuffled the total 400 data with a fixed seed in order to make sure the randomness, but could be stable reproduction.

2. We have known that the gradient of logistic regression from Problem 1 is:

$$\nabla E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n)$$

So the learning rate is set to be $\eta = 0.1$ $\eta_t = \eta \cdot \|\nabla E_{in}(\mathbf{w})\|$

And the total iterated steps is set to be ????

3. The results of the last 10 steps are shown as follow:

TODO: 1. last 10 steps' output 2. loss against iteration, gradient norm against iteration, stepsize against iteration 3. AUC, ROC curve

The code and the method to run the code are all in the folder 'code'.

Fixed learning rate gradient descent:

- 1: Initialize the weights at time step $t = 0$ to $\mathbf{w}(0)$.
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: Compute the gradient $\mathbf{g}_t = \nabla E_{\text{in}}(\mathbf{w}(t))$.
- 4: Set the direction to move, $\mathbf{v}_t = -\mathbf{g}_t$.
- 5: Update the weights: $\mathbf{w}(t + 1) = \mathbf{w}(t) + \eta \mathbf{v}_t$.
- 6: Iterate to the next step until it is time to stop.
- 7: Return the final weights.

Figure 3: The last 10 steps' output