

Machine Learning and Computational Statistics
Theodoros Georgiopoulos
Homework 8

Exercises for Unit 8: Logistic Regression classifier

Solution for exercise 1

We have a dataset $Y = \{(y_i, \mathbf{x}_i, i = 1, \dots, N)\}$ where $y_i \in \{0, 1\}$ is the class label for vector $\mathbf{x}_i \in R^l$. Lets extract the gradient descent logistic regression classifier. Because we have binary classification, the model is

$$\ln \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} = \boldsymbol{\theta}^T \mathbf{x}$$

which becomes

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{\boldsymbol{\theta}^T \mathbf{x}}} = \sigma(\boldsymbol{\theta}^T \mathbf{x})$$

using $P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x})$. The aim is to determine the $\boldsymbol{\theta}$ that maximizes $P(y_i = 1|\mathbf{x}_i)$ for all \mathbf{x}_i with $y_i = 1$ and maximizes $P(y_i = 0|\mathbf{x}_i)$ for all \mathbf{x}_i with $y_i = 0$. So, we want to maximize the likelihood

$$\prod_{\mathbf{x}_i: y_i=1} P(y_i = 1|\mathbf{x}_i) \prod_{\mathbf{x}_i: y_i=0} P(y_i = 0|\mathbf{x}_i).$$

We can make it simpler and maximize

$$\prod_{i=1}^N P(y_i = 1|\mathbf{x}_i)^{y_i} P(y_i = 0|\mathbf{x}_i)^{1-y_i} = \prod_{i=1}^N \sigma(\boldsymbol{\theta}^T \mathbf{x}_i)^{y_i} (1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}_i))^{1-y_i}.$$

Because we will use gradient descent, we want to minimize the negative log likelihood, which is

$$\begin{aligned} L(\boldsymbol{\theta}) &= - \sum_{i=1}^N \left(\sigma(\boldsymbol{\theta}^T \mathbf{x}_i)^{y_i} + (1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}_i))^{1-y_i} \right) \\ &= - \sum_{i=1}^N \left(y_i \ln(\sigma(\boldsymbol{\theta}^T \mathbf{x}_i)) + (1 - y_i) \ln(1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}_i)) \right). \end{aligned}$$

Taking the gradient with respect to θ and using $\frac{\partial \sigma(\theta^T \mathbf{x}_i)}{\partial \theta} = \sigma(\theta^T \mathbf{x}_i)(1 - \sigma(\theta^T \mathbf{x}_i))$ gives

$$\begin{aligned}
\nabla_{\theta} L(\theta) &= - \sum_{i=1}^N \left(\frac{y_i}{\sigma(\theta^T \mathbf{x}_i)} \sigma(\theta^T \mathbf{x}_i)(1 - \sigma(\theta^T \mathbf{x}_i)) \mathbf{x}_i - \frac{1 - y_i}{1 - \sigma(\theta^T \mathbf{x}_i)} \sigma(\theta^T \mathbf{x}_i)(1 - \sigma(\theta^T \mathbf{x}_i)) \mathbf{x}_i \right) \\
&= - \sum_{i=1}^N \left(y_i(1 - \sigma(\theta^T \mathbf{x}_i)) \mathbf{x}_i - (1 - y_i) \sigma(\theta^T \mathbf{x}_i) \mathbf{x}_i \right) \\
&= - \sum_{i=1}^N \left(y_i(1 - \sigma(\theta^T \mathbf{x}_i)) - (1 - y_i) \sigma(\theta^T \mathbf{x}_i) \right) \mathbf{x}_i \\
&= - \sum_{i=1}^N (y_i - \sigma(\theta^T \mathbf{x}_i)) \mathbf{x}_i \\
&= \mathbf{X}^T (\sigma(\theta^T \mathbf{x}) - \mathbf{y})
\end{aligned}$$

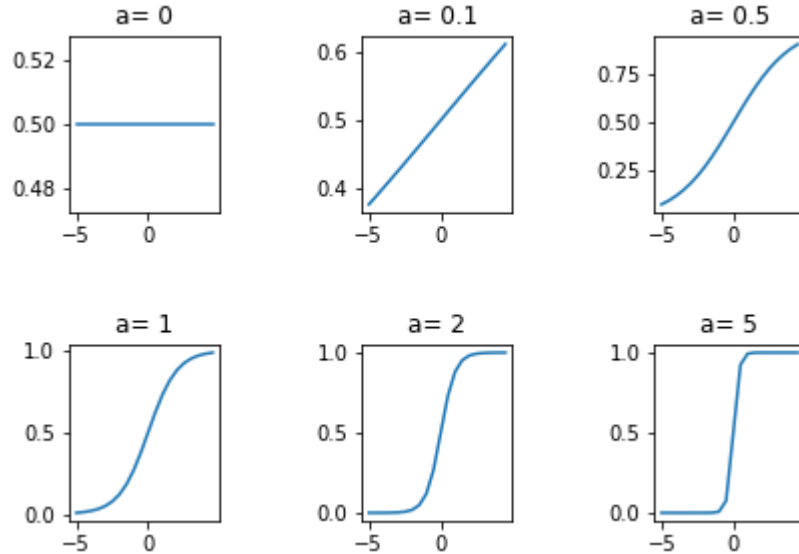
where $\mathbf{X}^T = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, $\mathbf{y} = [y_1, y_2, \dots, y_N]$, $\sigma(\theta^T \mathbf{x}) = [\sigma(\theta^T \mathbf{x}_1), \sigma(\theta^T \mathbf{x}_2), \dots, \sigma(\theta^T \mathbf{x}_N)]$. So, for the steps of gradient descent we have

$$\theta_i = \theta_{i-1} - \mathbf{X}^T (\sigma(\theta^T \mathbf{x}) - \mathbf{y}).$$

Solution for exercise 2

We have a dataset $Y = \{(y_i, \mathbf{x}'_i, i = 1, \dots, N)\}$ where $y_i \in \{0, 1\}$ is the class label for vector $\mathbf{x}'_i \in R^l$. The y and \mathbf{x}' are related such as: $y = f(\theta^T \mathbf{x}' + \theta_0)$ and $f(z) = 1/(1 + \exp(-az))$.

(a) We plot the function $f(z)$ for various values of the a .



(b) The sum of error squares criterion is

$$J(\theta) = \sum_{i=1}^N \left(y_i - f(\theta^T \mathbf{x}'_i) \right)^2$$

and the gradient with respect to $\boldsymbol{\theta}$ is

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = 2a \sum_{i=1}^N \left(y_i - f(\boldsymbol{\theta}^T \mathbf{x}'_i) \right) \left(f(\boldsymbol{\theta}^T \mathbf{x}'_i) (1 - f(\boldsymbol{\theta}^T \mathbf{x}'_i)) \right) \mathbf{x}'_i$$

and we can use it for the gradient step

$$\boldsymbol{\theta}_i = \boldsymbol{\theta}_{i-1} - \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_{i-1}).$$

- (c) We can see that for clear 1 response of the model, we need $\exp(-a\boldsymbol{\theta}^T \mathbf{x}'_i) = 0$ which is not possible. Moreover, even if $\exp(-a\boldsymbol{\theta}^T \mathbf{x}'_i) \rightarrow +\infty$, the model can't response a clear 0.
- (d) For a given \mathbf{x} , if $f(\boldsymbol{\theta}^T \mathbf{x}'_i) > 0.5$, we classify it to the class $y = 1$ and for \mathbf{x} , if $f(\boldsymbol{\theta}^T \mathbf{x}'_i) < 0.5$, we classify it to the class $y = 0$.
- (e) A way for leading the model responses very close to 1 (for class 1 vectors) or 0 (for class 0 vectors) could be to increase the parameter a to approximate step function.