

Optimizing model likelihood using gradient descent

0.1 Linear regression

Input Data:

$$\text{Independent: } \vec{x}^{(i)} \tag{1}$$

$$\text{Dependant: } y^{(i)} \tag{2}$$

Parameters to learn

$$\text{Weights: } \vec{w} \tag{3}$$

$$\text{Bias: } b \tag{4}$$

Model

$$\hat{y}^{(i)} \sim \mathcal{N}(\vec{w}^T \vec{x}^{(i)} + b, \sigma) \tag{5}$$

Error to minimize is negative log likelihood (assuming Gaussian)

$$E^{(i)} = -\log p(y^{(i)}|\vec{x}^{(i)}) \quad (6)$$

$$E^{(i)} = -\log \mathcal{N}(y^{(i)} | \vec{w}^T \vec{x}^{(i)} + b, \sigma) \quad (7)$$

$$E = -\sum_{i=1}^n \log p(y^{(i)}|\vec{x}^{(i)}) \quad (8)$$

$$\vec{w}, b = \arg \min_{\vec{w}, b} -\sum_{i=1}^n \log \mathcal{N}(y^{(i)} | \vec{w}^T \vec{x}^{(i)} + b, \sigma) \quad (9)$$

$$\vec{w}, b = \arg \min_{\vec{w}, b} -\sum_{i=1}^n \log \frac{1}{\sqrt{2\sigma^2\pi}} \exp -\frac{(y^{(i)} - \vec{w}^T \vec{x}^{(i)} - b)^2}{\sigma^2} \quad (10)$$

$$\vec{w}, b = \arg \min_{\vec{w}, b} \sum_{i=1}^n (y^{(i)} - \vec{w}^T \vec{x}^{(i)} - b)^2 \quad (11)$$

$$\frac{\partial E}{\partial \vec{w}} = \frac{2}{n} \sum_{i=1}^n (\vec{w}^T \vec{x}^{(i)} + b - y^{(i)}) \vec{x}^{(i)} \quad (12)$$

$$\frac{\partial E}{\partial \vec{w}} = \frac{2}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)}) \vec{x}^{(i)} \quad (13)$$

$$\frac{\partial E}{\partial b} = \frac{2}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)}) \quad (14)$$

0.2 Logistic regression

Input Data:

$$\text{Features: } \vec{x}^{(i)} \quad (15)$$

$$\text{Labels: } y^{(i)} \in \{1, 0\} \quad (16)$$

Parameters to learn

$$\text{Weights: } \vec{w} \quad (17)$$

$$\text{Bias: } b \quad (18)$$

Derived Features

$$z^{(i)} = \vec{w}^T \vec{x}^{(i)} + b \quad (19)$$

$$p^{(i)} = \sigma(z^{(i)}) \quad (20)$$

$$p^{(i)} = p(y^{(i)} = 1 | \vec{x}^{(i)}) = \sigma(z^{(i)}) \quad (21)$$

$$\sigma(a) = \frac{1}{1 + e^{-a}} \quad (22)$$

Likelihood of y comes from a Bernouli process parameterized by $p^{(i)}$

$$p(y^{(i)} | \vec{x}^{(i)}) = (p^{(i)})^{y^{(i)}} (1 - p^{(i)})^{(1 - y^{(i)})} \quad (23)$$

$$p(y | \vec{x}) = \prod_{i=1}^n p(y^{(i)} | \vec{x}^{(i)}) \quad (24)$$

Error to minimize is negative log likelihood

$$E^{(i)} = -(y^{(i)} \log p^{(i)} + (1 - y^{(i)}) \log(1 - p^{(i)})) \quad (25)$$

Use the chain rule to optimize \vec{w} and b

$$\frac{\partial E^{(i)}}{\partial \vec{w}} = \frac{\partial E^{(i)}}{\partial p^{(i)}} \frac{\partial p^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial \vec{w}} \quad (26)$$

$$\frac{\partial E^{(i)}}{\partial p^{(i)}} = -\frac{y^{(i)}}{p^{(i)}} + \frac{(1 - y^{(i)})}{(1 - p^{(i)})} \quad (27)$$

$$\frac{\partial p^{(i)}}{\partial z^{(i)}} = p^{(i)}(1 - p^{(i)}) = \sigma(z^{(i)})(1 - \sigma(z^{(i)})) \quad (28)$$

$$\frac{\partial z^{(i)}}{\partial \vec{w}} = \vec{x}^{(i)} \quad (29)$$

$$\frac{\partial E^{(i)}}{\partial \vec{w}} = (p^{(i)} - y^{(i)}) \vec{x}^{(i)} \quad (30)$$

$$\frac{\partial E^{(i)}}{\partial b} = \frac{\partial E^{(i)}}{\partial p^{(i)}} \frac{\partial p^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial b} \quad (31)$$

$$\frac{\partial z^{(i)}}{\partial b} = 1 \quad (32)$$

$$\frac{\partial E^{(i)}}{\partial b} = (p^{(i)} - y^{(i)}) \quad (33)$$

In gradient descent for batch B our update rule is

$$\vec{w}' = \vec{w} - \epsilon \frac{1}{|B|} \sum_{i \in B} \frac{\partial E^{(i)}}{\partial \vec{w}} \quad (34)$$

$$b' = b - \epsilon \frac{1}{|B|} \sum_{i \in B} \frac{\partial E^{(i)}}{\partial b} \quad (35)$$