

Weighted Least Square, Logistic Regression

Newton method

★ Probabilistic interpretation

⇒ Let us assume that the target variables and the inputs are related via the equation:

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

} Error term that captures
Unmodeled effects or
random noise

} $y^{(i)}, x^{(i)}, \epsilon^{(i)}$ are random variables

⇒ Let us further assume that the $\epsilon^{(i)}$ are distributed IID.

→ (Independent and Identically
Distributed)

$$\Rightarrow \epsilon^{(i)} \sim N(0, \sigma^2)$$

$$P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

or $Y^{(i)} | X^{(i)}; \theta \sim N(\theta^T X^{(i)}, \sigma^2)$

Let

$$X = \begin{bmatrix} X^{(1)T} \\ X^{(2)T} \\ \vdots \\ X^{(m)T} \end{bmatrix}$$

$$Y = \begin{bmatrix} Y^{(1)} \\ Y^{(2)} \\ \vdots \\ Y^{(m)} \end{bmatrix}$$

\Rightarrow The probability of the data is given by $P(Y|X; \theta)$

\Rightarrow This quantity is typically viewed as a function of Y , for a fixed value of θ .

\Rightarrow When we wish to explicitly view this as a function of θ , we will instead call it the **likelihood** function:

$$L(\theta) = L(\theta; X, Y) = P(Y|X; \theta)$$

$$\Rightarrow L(\theta) = \prod_{i=1}^m P(Y^{(i)} | X^{(i)}; \theta) \quad \left\{ \varepsilon^{(i)} \text{ are independent} \right\}$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(Y^{(i)} - \theta^T X^{(i)})^2}{2\sigma^2}\right)$$

⇒ We should choose θ to maximize $L(\theta)$.

⇒ Instead of maximizing $L(\theta)$, we can also minimize any strictly increasing function of $L(\theta)$.

⇒ Let $l(\theta) = \log L(\theta)$ {log likelihood}

$$l(\theta) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$= m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

⇒ Hence, maximizing $l(\theta)$ gives the same answer as minimizing:

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

* Locally weighted linear regression

⇒ It assumes that there are sufficient training data, and makes the choice of features less critical.

⇒ An locally weighted linear regression algorithm does the following:

① Fit θ to minimize $\sum_i \omega^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$

② Output $\theta^T x$

A fairly standard choice for the weights is

$$\omega^{(i)} = \exp\left(-\frac{\|x^{(i)} - x\|^2}{2\tau^2}\right)$$

} Controls how quickly the weight of a training example falls off with distance of its $x^{(i)}$ from the query pair x .

} τ is called the **bandwidth parameter**

* Classification and Logistic regression

Given $x^{(i)}$, the corresponding $y^{(i)}$ is called the label for the training example.

We could approach the classification problem ignoring the fact that y is discrete-valued, and use old linear regression algorithm to try to predict y given x .

However, it is easy to construct examples where this method performs very poorly.

$$\Rightarrow \text{Let } h_0(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

where, $g(z) = \frac{1}{1 + e^{-z}}$ is logistic function

Given the logistic regression model how do we fit θ for it?

Lets endow our classification model with a set of probabilistic assumptions, and then fit the parameters via maximum likelihood.

Let us assume that:

$$P(y=1 | x; \theta) = h_0(x)$$

$$P(y=0 | x; \theta) = 1 - h_0(x)$$

$$\Rightarrow P(y|x; \theta) = (h_\theta(x))^y (1 - h_\theta(x))^{1-y}$$

\Rightarrow Assuming the m training examples were generated independently, we can then write down the likelihood of the parameter as:

$$L(\theta) = P(Y|X; \theta)$$

$$= \prod_{i=1}^m P(y^{(i)}|x^{(i)}; \theta)$$

$$= \prod_{i=1}^m (h_\theta(x^{(i)}))^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}}$$

\Rightarrow As before, it will be easier to maximize the log likelihood:

$$l(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))$$

\Rightarrow To maximize the log likelihood we can use gradient ascent.

$$\theta := \theta + \alpha \nabla_{\theta} l(\theta)$$

$$\boxed{\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}}$$

{ This is a property of much bigger class of problems called generalized linear model }

* Another algorithm for maximizing $\ell(\theta)$

{ Newton's method }

\Rightarrow Let $f: \mathbb{R} \rightarrow \mathbb{R}$, and we wish to find a value of θ so that $f(\theta) = 0 \cdot \forall \theta \in \mathbb{R}$

\Rightarrow Newton's method performs the following update:

so

$$\theta := \theta - \frac{f(\theta)}{f'(\theta)}$$

\Rightarrow This method has a natural interpretation in which we can think of it as approximating the function f via a linear function

\rightarrow That is tangent to f at the current guess θ

\rightarrow Solving for where that linear function equals to zero

\rightarrow And letting the next guess for θ be where that linear function is zero.

\Rightarrow Suppose we want to maximize some function ℓ .

\Rightarrow Maxima of ℓ corresponds to points where its first derivative $\ell'(\theta)$ is zero.

⇒ So by letting $f(\theta) = l'(\theta)$, we can use the same algorithm to maximize l , and we obtain update rule:

$$\theta := \theta - \frac{l'(\theta)}{l''(\theta)}$$

⇒ Lastly, in our logistic regression setting, θ is vector-valued, so we need to generalize Newton's method to this setting.

⇒ The generalization of Newton's method to this multidimensional setting is given by

$$\theta := \theta - H^{-1} \nabla_{\theta} l(\theta)$$

⇒ When Newton's method is applied to maximize the logistic regression log likelihood function $l(\theta)$, the resulting method is also called **Fisher Scoring**.

~~Let $\nabla_{\theta} l(\theta) = f(\theta) = 0$ {for maximizing $l(\theta)$ }~~

$$f(\theta + \Delta\theta) \approx f(\theta) + \frac{\delta f}{\delta \theta} \Delta\theta = 0$$

$$= \nabla_{\theta} f(\theta) + H(\theta) \Delta\theta = 0$$

$$\Delta\theta = -H^{-1}(\theta) \nabla_{\theta} l(\theta)$$

So $\theta := \theta - H^{-1}(\theta) \nabla_{\theta} l(\theta)$

Hessian matrix

* The perceptron learning algorithm

⇒ Consider modifying the logistic regression method to "force" it to output values that are either 0 or 1.

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

⇒ If we then let $h_\theta(x) = g(\theta^T x)$ and if we use the update rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

⇒ We have perceptron learning algorithm.

⇒ In the 1960s, this perceptron was argued to be a rough model for how individual neurons in the brain work.

⇒ It is difficult to endow the perceptron's predictions with meaningful probabilistic interpretation, or derived the perceptron as a maximum likelihood estimation algorithm.