

Logistics Regression

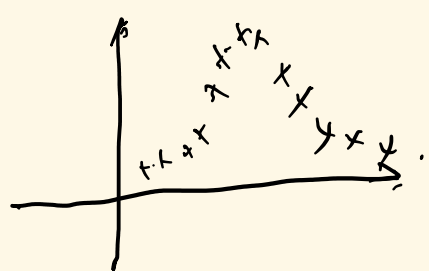
Monday, 19 March 2018

08:33

underfitting / overfitting.

Parametric learning algorithm / non parametric learning.

Locally weight regression: (LWSS)



LWR: Fit θ to minimize

$$\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

$$\text{where } w^{(i)} = \exp\left(-\frac{(x^{(i)} - \bar{x})^2}{2\tau^2}\right)$$

τ : bandwidth parametric.

Short coming: training every time when querying.

* Fix: Andrew Moore: KD Tree.

Probabilistic Interpretation:

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

$$\epsilon^{(i)} = \text{error}, \quad \epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

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

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

$$\Rightarrow y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)$$

CLT permit error to be normal distribution.

assume: $\epsilon^{(i)}$ are IID (Independently and Identically Distributed)

$$L(\theta) = p(\vec{y} | X; \theta): \text{likelihood.}$$

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

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

Maximize likelihood:

choose θ to maximize $L(\theta)$

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^m \log\left[\frac{1}{\sqrt{2\pi}\sigma} \exp(\dots)\right]$$

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

to maximize $\ell(\theta)$, we need minimize $\sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$

Variance σ will take no effect on θ .

Classification:

Linear regression works not very well on classification.

Choose: $h_{\theta}(x) \in [0, 1]$.

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad // \quad g(x) = \frac{1}{1 + e^{-x}}$$

logistic / sigmoid function.

$$p(y=1 | x; \theta) = h_{\theta}(x)$$

$$p(y=0 | x; \theta) = 1 - h_{\theta}(x)$$

$$p(y | x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y} \quad \text{more elegant way}$$

$$L(\theta) = p(\vec{y} | X; \theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

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

// Gradient ascent

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

$$\frac{\partial \ell}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Perceptron Algorithm:

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

$$h_{\theta}(x) = g(\theta^T x)$$

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