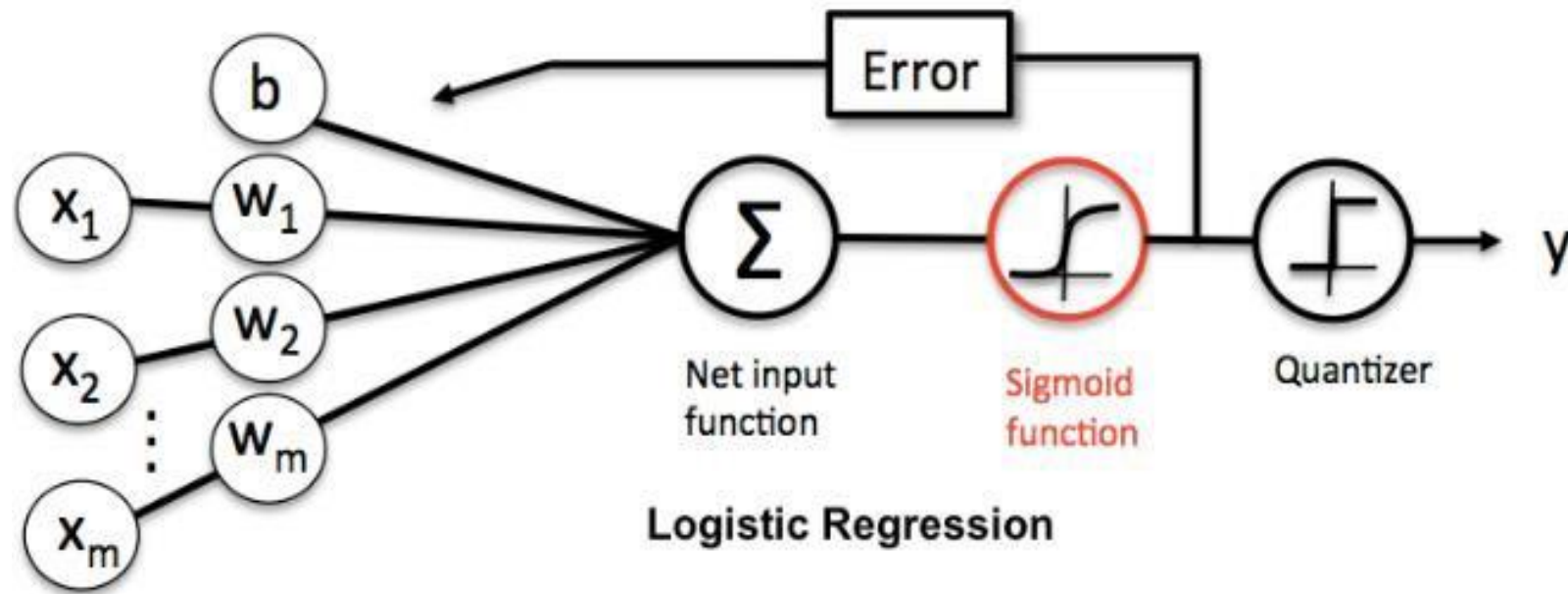




The Logistic Regressor

Alexander G. Ororbia II
Introduction to Machine Learning
CSCI-635
10/4/2023

Logistic Regressor Architecture



(Representation!)

Where Does the **Form** Come From?

- Logistic regression hypothesis representation

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)}}$$

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features $[X_1, \dots, X_n]^T$
 - Y is Boolean
 - Assume all X_i are **conditionally independent** given Y
 - Model $P(X_i | Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i)$
 - Model $P(Y)$ as Bernoulli π

What is $P(Y | X_1, X_2, \dots, X_n)$?

(Representation!)



Logistic Regression

- Hypothesis representation
- **Cost function**
- Logistic regression with gradient descent
- Regularization
- Multi-class classification

Training set with m examples

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$$

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

How to choose parameters θ ?

Cost function for Linear Regression

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y)$$

$$\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$

Logistic Regression Cost Function

- $\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$



- $\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$

- If $y = 1$: $\text{Cost}(h_{\theta}(x), y) = -\log(h_{\theta}(x))$

- If $y = 0$: $\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x))$

Logistic Regression Cost Function

- $$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

**But, where does
it come from?**



- $$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

- If $y = 1$: $\text{Cost}(h_{\theta}(x), y) = -\log(h_{\theta}(x))$

- If $y = 0$: $\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x))$

Logistic Regression

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] \end{aligned}$$

Learning: fit parameter θ

$$\min_{\theta} J(\theta)$$

Prediction: given new x

$$\text{Output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Where Does the **Cost** Come From?

- Training set with m examples

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

- Maximum likelihood estimate for parameter θ

$$\begin{aligned}\theta_{\text{MLE}} &= \operatorname{argmax}_{\theta} P_{\theta} \left((x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}) \right) \\ &= \operatorname{argmax}_{\theta} \prod_{i=1}^m P_{\theta} \left((x^{(i)}, y^{(i)}) \right)\end{aligned}$$

- Maximum conditional likelihood estimate for parameter θ

Remember the Bernoulli Distribution?

$$P(\mathbf{x} = 1) = \phi$$

$$P(\mathbf{x} = 0) = 1 - \phi$$

$$P(\mathbf{x} = x) = \phi^x (1 - \phi)^{1-x}$$

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}] = \phi$$

$$\text{Var}_{\mathbf{x}}(\mathbf{x}) = \phi(1 - \phi)$$

- **Goal:** choose θ to maximize conditional likelihood of training data

- $P_{\theta}(Y = 1|X = x) = h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$

- $P_{\theta}(Y = 0|X = x) = 1 - h_{\theta}(x) = \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}$

- **Training data** $D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

- **Data likelihood** $= \prod_{i=1}^m P_{\theta}((x^{(i)}, y^{(i)}))$

- **Data conditional likelihood** $= \prod_{i=1}^m P_{\theta}(y^{(i)}|x^{(i)})$

$$\theta_{\text{MCLE}} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P_{\theta}(y^{(i)}|x^{(i)})$$

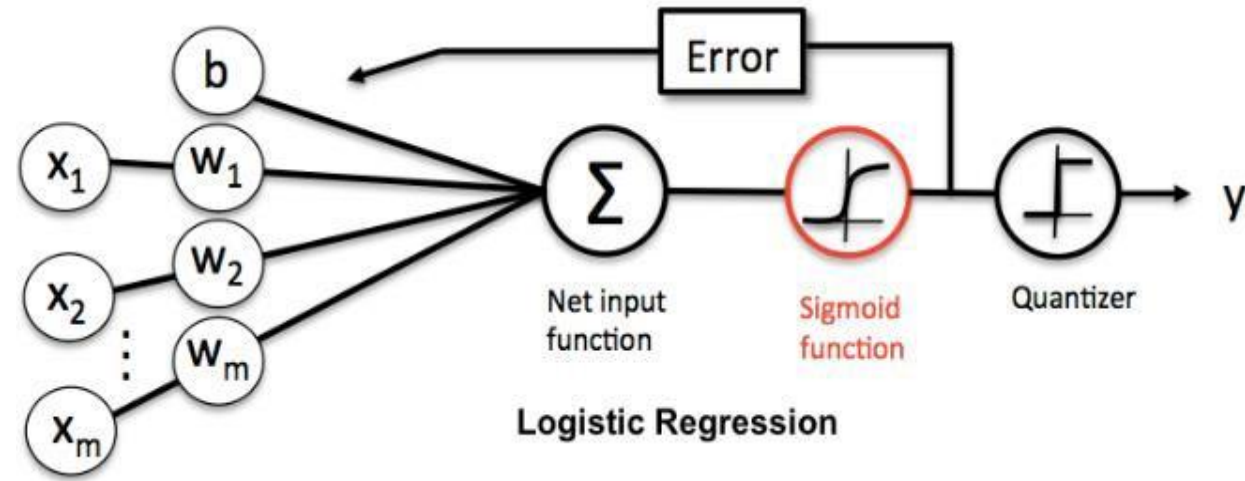
Expressing Conditional Log-Likelihood

$$\begin{aligned} L(\theta) &= \log \prod_{i=1}^m P_{\theta}(y^{(i)} | x^{(i)}) = \sum_{i=1}^m \log P_{\theta}(y^{(i)} | x^{(i)}) \\ &= \sum_{i=1}^m y^{(i)} \log P_{\theta}(y^{(i)} = 1 | x^{(i)}) + (1 - y^{(i)}) \log P_{\theta}(y^{(i)} = 0 | x^{(i)}) \\ &= \sum_{i=1}^m y^{(i)} \log (h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \end{aligned}$$

The logarithm is
your *friend*!

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

Logistic Regression



$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

Bernoulli log likelihood!

Learning: fit parameter θ

$$\min_{\theta} J(\theta)$$

Prediction: given new x

$$\text{Output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

(Evaluation!)

Gradient Descent

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

Goal: $\min_{\theta} J(\theta)$

Good news: Convex function!

Bad news: No analytical solution

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

What do we still need?

Derive gradient of $J(\theta_j)$ with respect to each θ_j

(Optimization!)

Questions?

Deep robots!

Deep questions?!

