

Logistic Regression: Classification

Email: Spam/Not?

Online Transactions: Fraudulent Yes/No

Tumor: Malign / Benign?

$y \in \{0, 1\}$

0: "Negative Class"
1: "Positive Class"

absence of presence of

multiclass $y \in \{0, 1, 2, 3\}$

binary

threshold classifier output $h_\theta(x)$ at 0.5:

If $h_\theta(x) \geq 0.5$, " $y = 1$ "

If $h_\theta(x) < 0.5$, " $y = 0$ "

Classification: $y = 0$ or 1

$h_\theta(x)$ can be > 1 or < 0

Logistic regression: $0 \leq h_\theta(x) \leq 1$

Classification

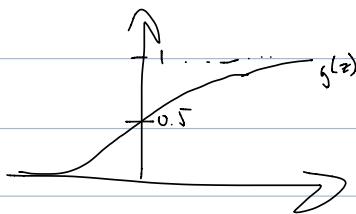


$$W_{int} + b \leq h_\theta(x) \leq 1$$

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

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

$$g(z) = \frac{1}{1 + e^{-z}} \rightarrow \text{sigmoid/logistic function}$$



Interpretation of Hypothesis Output

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x

$$h_{\theta}(x) = P(y=1 \mid x; \theta) \quad \text{"probability that } y=1 \text{, given } x, \text{ parameterized by } \theta"$$

$y = 0 \text{ or } 1$

$$P(y=0 \mid x; \theta) + P(y=1 \mid x; \theta) = 1$$

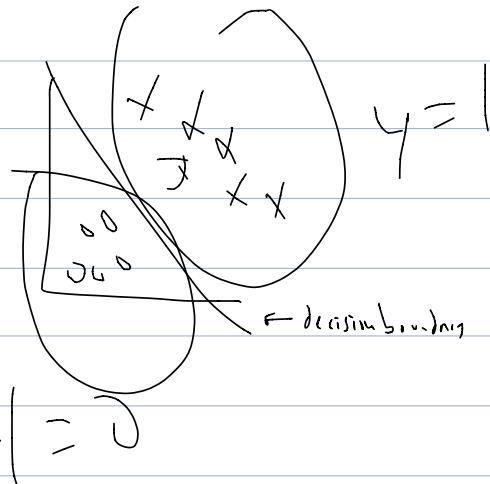
$$P(y=0 \mid x; \theta) = 1 - P(y=1 \mid x; \theta)$$

Decision Boundary

$$g(z) \geq 0.5 \text{ when } z \geq 0$$

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5 \text{ when } \theta^T x \geq 0$$

$$\theta = \begin{bmatrix} -3 \\ 1 \end{bmatrix} \rightarrow \text{with "y=1" if } \underbrace{-3 + x_1 + x_2}_{\theta^T x} \geq 0 \\ x_1 + x_2 \geq 3$$



Logistic Regression : Cost Function

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$

in n examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad x_0 = 1, y \in \{0, 1\}$$

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

How to choose parameters θ ?

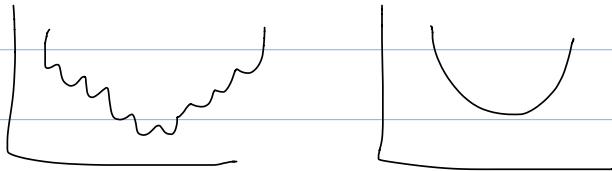
Linear regression: $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$

Logistic

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

"non-convex"

"convex"



Logistic Regression Cost Function

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

Simplified Cost Function & gradient descent

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

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

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

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

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

To fit parameters θ :

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

To make a prediction given x :

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

Gradient descent

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

Repeat {

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

}

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

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

Repeat {

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

}

Algorithm looks identical to linear regression!

Advanced Optimization

• Gradient descent

• Conjugate gradient

• BFGS

• L-BFGS

Multiclass Classification

ex: Email filtering/tags: Work, Family, Hobby

(x2): medical diagnosis: not ill, cold, flu

(x3): Weather: sunny, cloudy, rain, snow

$$y = \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix}$$

One vs all (One vs rest)

$$h_{\theta}^{(i)}(x) = P(y=i|x; \theta) \quad (i=1,2,3)$$

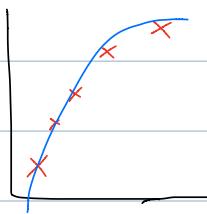


The Problem of Overfitting



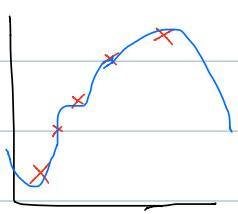
$$\theta_0 + \theta_1 x$$

"underfit" "high bias"



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"just right"



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

"overfit" "high variance"

Addressing Overfitting

options:

1. Reduce # of features

- manually select

- Model selection algorithm

2. Regularization

- keep all features, but reduce magnitude/values of parameters θ_j

Regularization Cost Function

Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$

- "Simple" hypothesis

- Loss function to overfitting

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

regulation parameter

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



Regularized Linear Regression

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$
$$\min_{\theta} J(\theta)$$

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\left. \begin{array}{l} \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \\ \quad \left(\frac{\partial}{\partial \theta_j} J(\theta) \right) \text{ regularization} \end{array} \right\}$$

$$\theta_j := \underbrace{\theta_j (1 - \alpha \frac{\lambda}{m})}_{< 1} - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\theta = (x^T x + \lambda \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix})^{-1} x^T y$$

Regularized Logistic Regression

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

}

$$\frac{\partial J(\theta)}{\partial \theta_j}$$

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