



Classification Algorithms

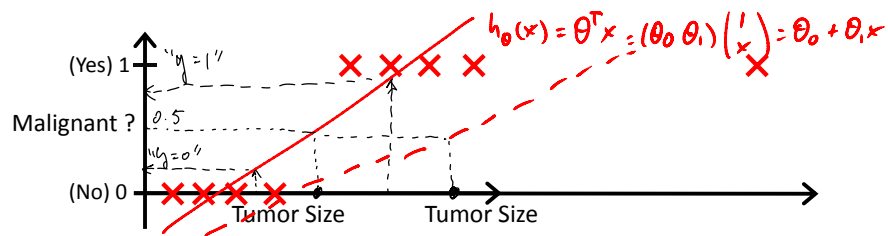
Logistic Regression

Classification

Email: Spam / Not Spam?
Online Transactions: Fraudulent (Yes / No)?
Tumor: Malignant / Benign ?

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

0: "Negative Class" (e.g., benign tumor)
1: "Positive Class" (e.g., malignant tumor)



Threshold classifier output $h_{\theta}(x)$ at 0.5:

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

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

because $y \in \{0, 1\}$
we want
 $0 \leq h_{\theta}(x) \leq 1$

Classification: $y = 0$ or 1

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

Logistic Regression: $0 \leq h_{\theta}(x) \leq 1$

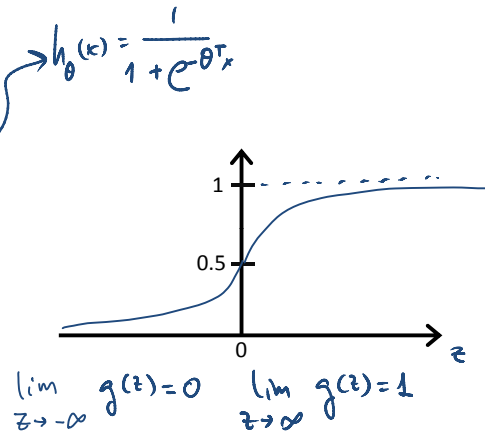
Logistic Regression Model

Want $0 \leq h_{\theta}(x) \leq 1$

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

$$g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid function
Logistic function



Interpretation of Hypothesis Output

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

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$$h_{\theta}(x) = 0.7$$

Tell patient that 70% chance of tumor being malignant

$h_{\theta}(x) = P(y=1|x; \theta) \rightarrow$ "probability that $y = 1$, given x , parameterized by θ "

because $g \in \{0, 1\} \rightarrow \begin{cases} P(y=0|x; \theta) + P(y=1|x; \theta) = 1 \\ P(y=0|x; \theta) = 1 - P(y=1|x; \theta) \end{cases}$

Logistic regression

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

$$g(z) = \frac{1}{1 + e^{-z}}$$

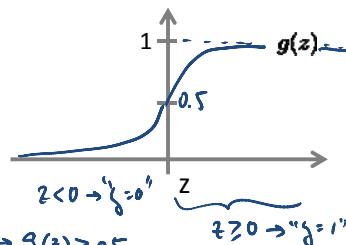
Suppose predict " $y = 1$ " if $h_{\theta}(x) \geq 0.5 \rightarrow g(z) \geq 0.5$

$$\hookrightarrow z \geq 0 \rightarrow \theta^T x \geq 0$$

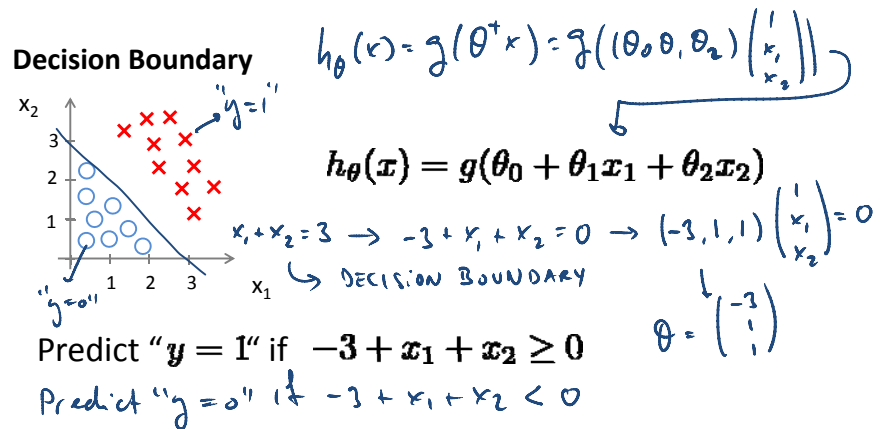
predict " $y = 0$ " if $h_{\theta}(x) < 0.5 \rightarrow g(z) < 0.5$

$$\hookrightarrow z < 0 \rightarrow \theta^T x < 0$$

DECISION BOUNDARY $\theta^T x = 0$



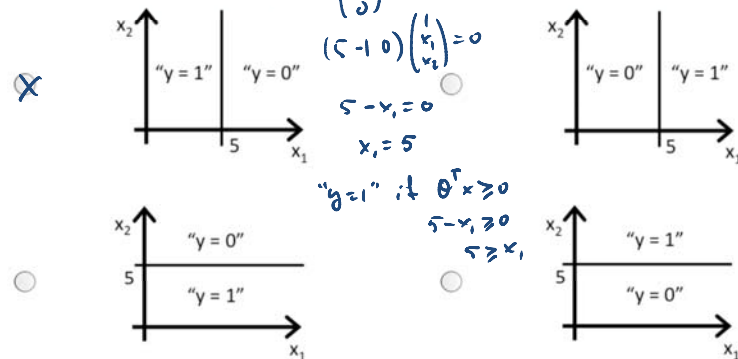
Decision Boundary



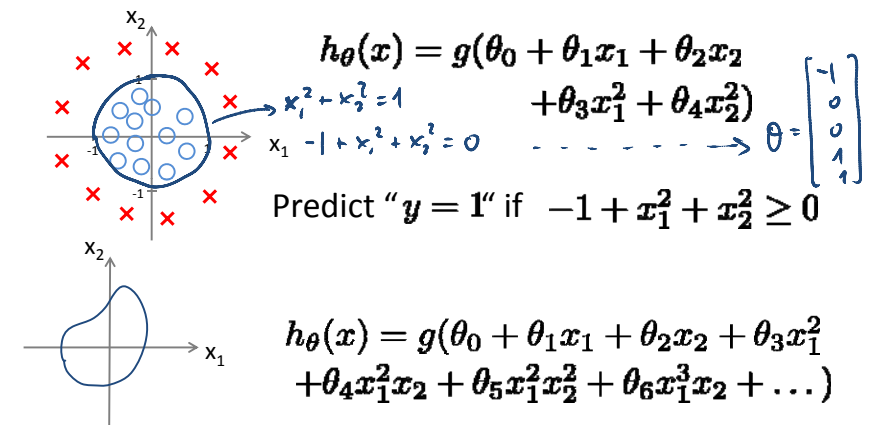
Predict " $y = 1$ " if $-3 + x_1 + x_2 \geq 0$

Predict " $y = 0$ " if $-3 + x_1 + x_2 < 0$

Suppose we use ~~linear~~^{logistic} regression with two characteristics x_1 y x_2 and get $\theta_0=5$, $\theta_1=-1$, $\theta_2=0$. Which of these figures shows the decision boundary?



Non-linear decision boundaries



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

m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix}$ $x_0 = 1, y \in \{0, 1\}$

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

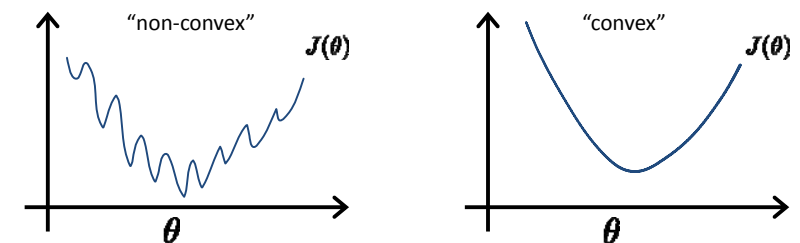
How to choose parameters θ ?

Cost function

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

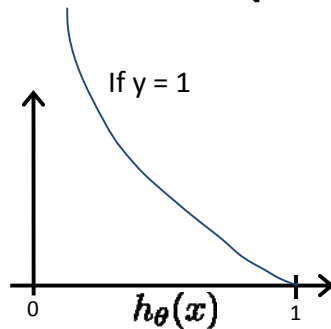
Cost($h_{\theta}(x^{(i)}), y^{(i)}$) = $\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$\frac{1}{1 + e^{-\theta^T x}}$



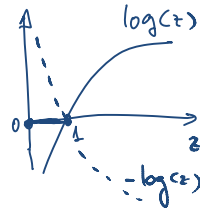
Logistic regression cost function

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



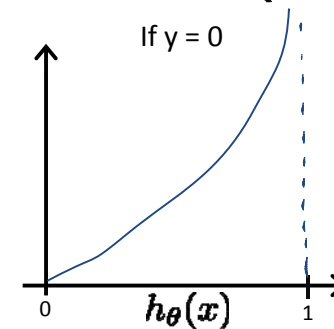
Cost = 0 if $y = 1, h_{\theta}(x) = 1$
 But as $h_{\theta}(x) \rightarrow 0 \rightarrow h_{\theta}(x) = P(y=1|x;\theta)$
 $\text{Cost} \rightarrow \infty$

Captures intuition that if $h_{\theta}(x) = 0$,
 (predict $P(y=1|x;\theta) = 0$), but $y = 1$,
 we'll penalize learning algorithm by a very large cost.



Logistic regression cost function

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



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

this cost is known as the
 $-\log$ likelihood

Logistic regression cost function

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

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

Note: $y = 0$ or 1 always

Logistic regression cost function

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

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\log h_{\theta}(x^{(i)})}_{\text{cost with } y^{(i)}=1} + (1 - y^{(i)}) \underbrace{\log(1 - h_{\theta}(x^{(i)}))}_{\text{cost with } y^{(i)}=0} \right]$$

To fit parameters θ :

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

To make a prediction given new x :

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

$$\rightarrow \begin{cases} \text{if } h_{\theta}(x) \geq 0.5 \rightarrow \theta^T x \geq 0 \rightarrow y = 1 \\ \text{if } h_{\theta}(x) < 0.5 \rightarrow \theta^T x < 0 \rightarrow y = 0 \end{cases}$$

Gradient Descent

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

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

Repeat {

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

} (simultaneously update all θ_j)

Gradient Descent

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

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

Repeat {

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

(simultaneously update all θ_j)

in logistic regression $h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$
in linear regression $h_{\theta}(x^{(i)}) = \theta^T x^{(i)}$
for regularization

Algorithm looks identical to linear regression!

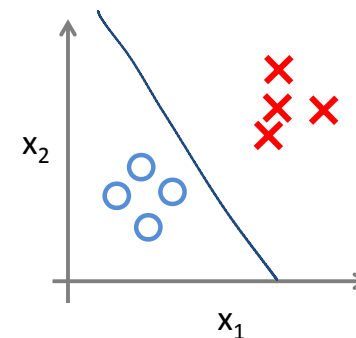
Multiclass classification

Email foldering/tagging: Work, Friends, Family, Hobby

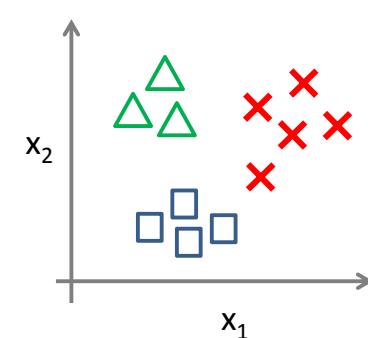
Medical diagrams: Not ill, Cold, Flu

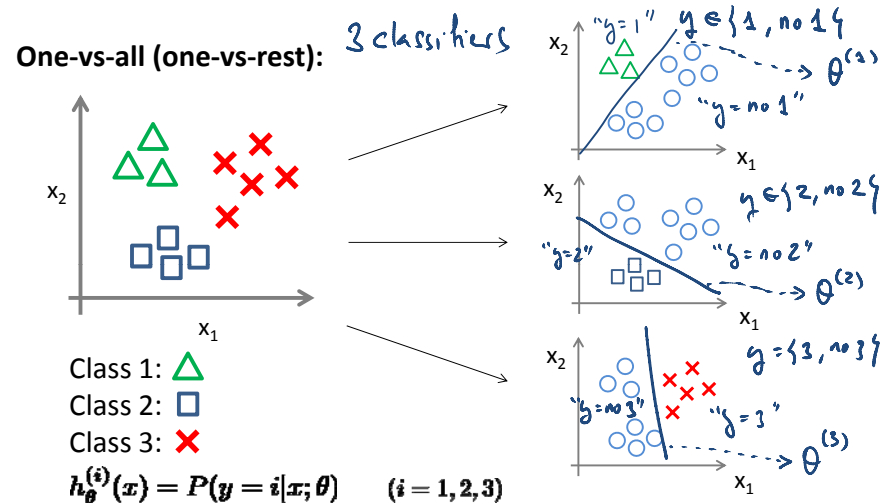
Weather: Sunny, Cloudy, Rain, Snow

Binary classification:



Multi-class classification:





One-vs-all

Train a logistic regression classifier $h_{\theta}^{(i)}(x)$ for each class i to predict the probability that $y = i$.
 $h_{\theta}^{(1)}(x) = h_{\theta^{(1)}}(x) = P(y=1|x;\theta^{(1)})$, $h_{\theta}^{(2)}(x) = P(y=2|x;\theta^{(2)}) \dots$
 On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(x) \rightarrow \max_i P(y=i|x;\theta^{(i)})$$

$$\hookrightarrow \max_i (\theta^{(i)})^T x$$

Suppose we have a classification problem with k classes. Using the 1-vs-all method, how many logistical sorters will we have to train?

- ☐ $k-1$
- ☐ k
- ☐ $k+1$
- ☐ Approximately $\log_2(k)$