# Introduction to ML
# Lecture 7: Logistic Regression
# (Linear Classifier)

Hamed Tabkhi

Department of Electrical and Computer Engineering,
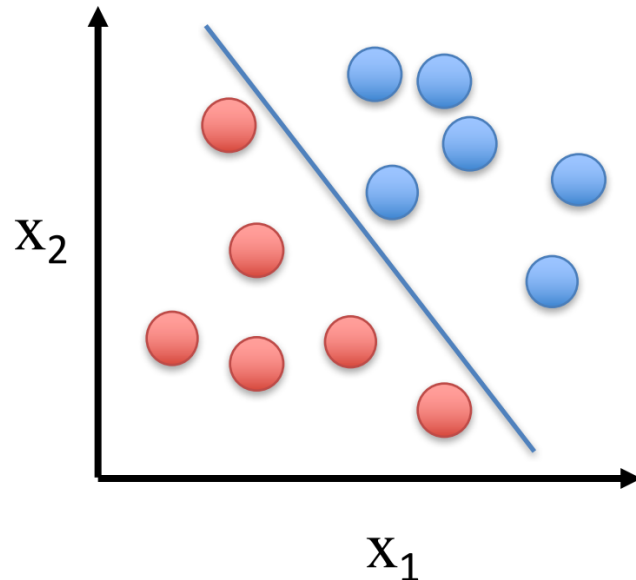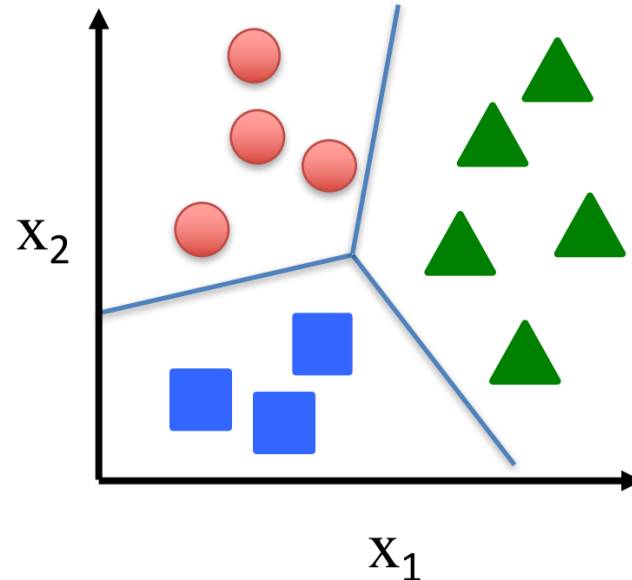University of North Carolina Charlotte (UNCC)
*htabkhiv@uncc.edu*

# Multi-classification (e.g. two explanatory variables)



Binary classification:

Multi-class classification:

Disease diagnosis:     healthy / cold / flu / pneumonia

Object classification:  desk / chair / monitor / bookcase

UNC CHARLOTTE

# Linear regression for classification

- We have discussed about regression
  - Output real value prediction

**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)

UNC CHARLOTTE

# Classification based on probability

- Instead of just predicting the class, give the probability of the instance being that class
  - i.e., learn $p(y \mid \boldsymbol{x})$

- Recall that:
$$0 \leq p(\text{event}) \leq 1$$
$$p(\text{event}) + p(\neg\text{event}) = 1$$

Not

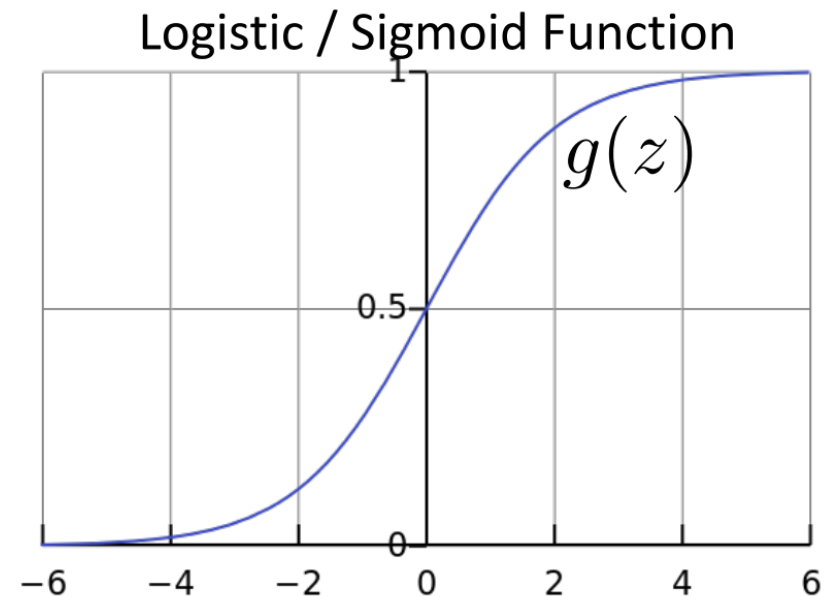Note: Although the name says "regression", but logistic regression is a classification approach

UNC CHARLOTTE

# Logistic regression

- Why sigmoid function

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

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

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

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

Logistic / Sigmoid Function

$g(z)$

UNC CHARLOTTE

# Logistic regression

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

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

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

- Assume a threshold and...

  – Predict y = 1 if $\boxed{h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5}$

  $$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x} \geq 0$$

  – Predict y = 0 if $\boxed{h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5}$

  $$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x} < 0$$

Logistic / Sigmoid Function



$g(z)$

UNC CHARLOTTE

# Logistic regression

- Interpretation of Hypothesis Output

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

$$p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$$

"probability that y = 1, given x, parameterized by $\theta$"

Example: Cancer diagnosis from tumor size

$$\boldsymbol{x} = \left[ \begin{array}{c} x_0 \\ x_1 \end{array} \right] = \left[ \begin{array}{c} 1 \\ \text{tumorSize} \end{array} \right]$$

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

→ Tell patient that 70% chance of tumor being malignant

This is a classification task, y takes on 1 or 0

Note that: $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) + p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1$
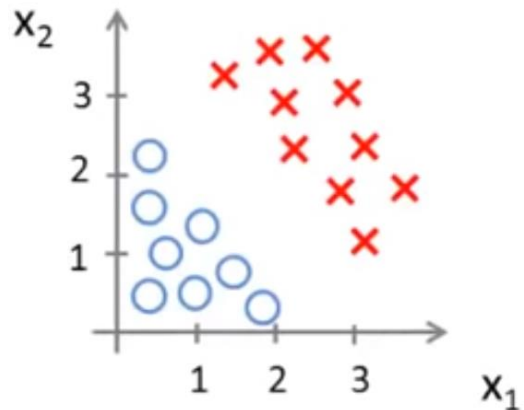
Therefore, $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1 - p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

UNC CHARLOTTE

# Logistic regression

- Example (let's see how the hypothesis is used to make predictions)

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

We haven't discussed how to fit the parameters $\boldsymbol{\theta}$
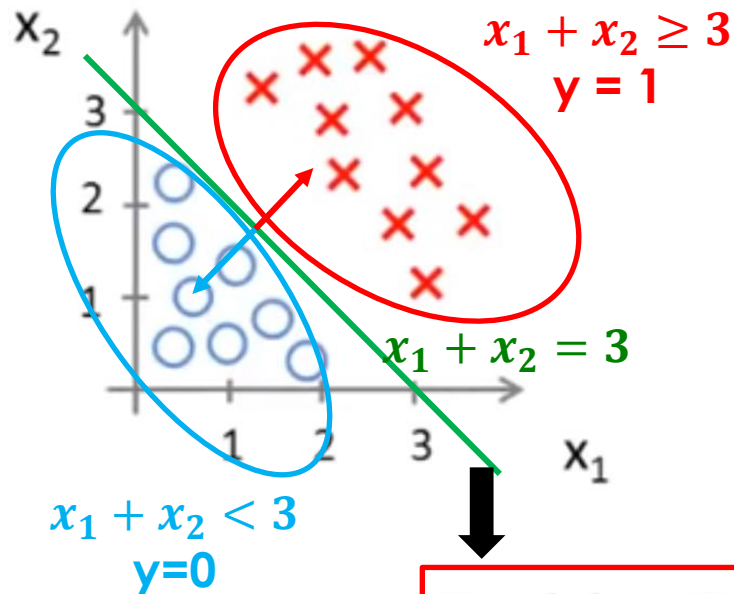
But, here let's assume we have $\boldsymbol{\theta} = [-3, 1, 1]^T$

$$h_\theta(x) = g(\underline{\theta_0} + \underline{\theta_1} x_1 + \underline{\theta_2} x_2)$$

-3  1   1

$$\boldsymbol{\theta}^T x = -3 + x_1 + x_2$$

UNC CHARLOTTE

# Logistic regression



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

$$x_1 + x_2 \geq 3$$

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

$x_1 + x_2 \geq 3$
y = 1

$x_1 + x_2 = 3$

$x_1 + x_2 < 3$
y=0

**Decision Boundary**

UNC CHARLOTTE

# Logistic regression – fit parameters

- Given $\left\{ \left( \boldsymbol{x}^{(1)}, y^{(1)} \right), \left( \boldsymbol{x}^{(2)}, y^{(2)} \right), \ldots, \left( \boldsymbol{x}^{(m)}, y^{(m)} \right) \right\}$   $m$ training samples

  where   $\boldsymbol{x}^{(i)} \in \mathbb{R}^n, \ y^{(i)} \in \{0, 1\}$

- Model:   $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left( \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x} \right)$

$$g(z) = \frac{1}{1 + e^{-z}}$$
    $\longrightarrow$     scales $\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}$ to [0, 1]

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \qquad \boldsymbol{x}^{\mathsf{T}} = \begin{bmatrix} 1 & x_1 & \ldots & x_n \end{bmatrix}$$

How to choose parameter $\boldsymbol{\theta}$ ?

UNC CHARLOTTE

# Logistic regression

**Cost function**

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

$\longrightarrow \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$

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

Logistic regression  $\qquad h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \dfrac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}}}$

Can we use the squared loss for logistic regression as in linear regression ?

UNC CHARLOTTE

# Logistic regression

**Logistic regression cost function**

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

This cost function is convex

UNC CHARLOTTE

# Logistic regression

$$\text{cost}\,(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} \boxed{-\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 1} \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 0 \end{cases}$$

$$0 \le h_{\boldsymbol{\theta}}(\boldsymbol{x}) \le 1$$

**Intuition behind the Objective**   If y = 1

- Cost = 0 if prediction is correct

If y = 1

cost

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$

0         1

Probability = 1 (100% that the label is 1, which is the ground truth), so let's don't impose any cost

e.g., probability = 0.8 (80% that the label is 1, which is a good prediction, but not perfect → add a small cost)
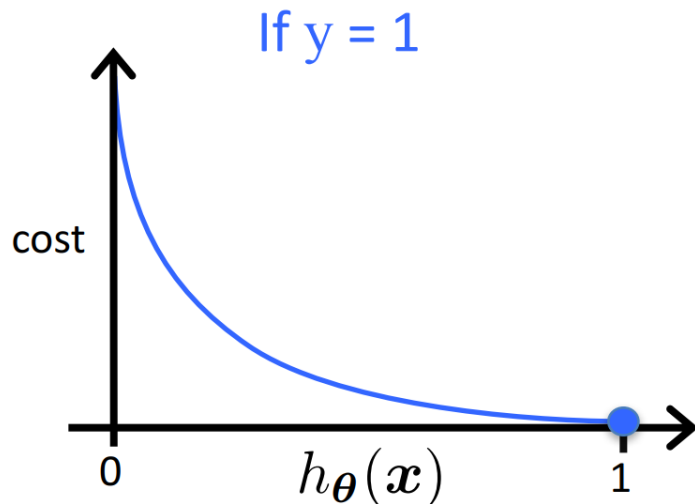
UNC CHARLOTTE

# Logistic regression

$$\text{cost}\,(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} \boxed{-\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 1} \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 0 \end{cases}$$

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

**Intuition behind the Objective**

If y = 1

- Cost = 0 if prediction is correct
- As $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties
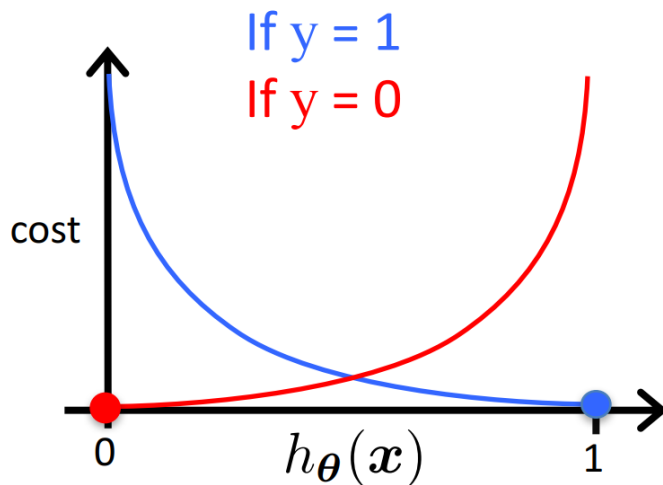  - e.g., predict $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0$, but y = 1

If y = 1



cost

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$

0          1

UNC CHARLOTTE

# Logistic regression

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

If y = 0

- Cost = 0 if prediction is correct

- As $(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties

If y = 1
If y = 0

cost

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$

0    1

UNC CHARLOTTE

# The cost function of logistic regression

**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

Compact form:

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

UNC CHARLOTTE

# Logistic regression

**Logistic regression cost function**

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

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

Find the parameters using Gradient descent

**Gradient Descent**

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

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

Repeat {

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

(simultaneously update all $\theta_j$)

}

UNC CHARLOTTE

# Logistic regression

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

- Logistic regression model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}\right)$$

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

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

You can do the math yourself, if you are interested!

$$
\begin{aligned}
g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\
&= \frac{1}{(1 + e^{-z})^2} \left(e^{-z}\right) \\
&= \frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right) \\
&= g(z)(1 - g(z)).
\end{aligned}
$$

See here:
https://towardsdatascience.com/derivative-of-the-sigmoid-function-536880cf918e

UNC CHARLOTTE

# Logistic regression

**Gradient Descent**

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

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

This looks IDENTICAL to linear regression!!!

Repeat {

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

(simultaneously update all $\theta_j$)

}

Another good reason of using Sigmoid function: mathematical convenient when computing the derivative

✦ However, the form of the model is very different:

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

UNC CHARLOTTE

# Gradient descent for Linear Regression

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \quad \boxed{h_\theta(x) = \theta^\top x}$$

}

# Gradient descent for Logistic Regression

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \quad \boxed{h_\theta(x) = \frac{1}{1 + e^{-\theta^\top x}}}$$

}

UNC CHARLOTTE

# Logistic regression

We can use gradient descent to learn parameter values, and hence compute the prediction for a new input.

To make a prediction given new $x$ :

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

= estimated probability that y = 1 on input x

UNC CHARLOTTE

# Logistic regression

- How to solve multi-class classification?

UNC CHARLOTTE
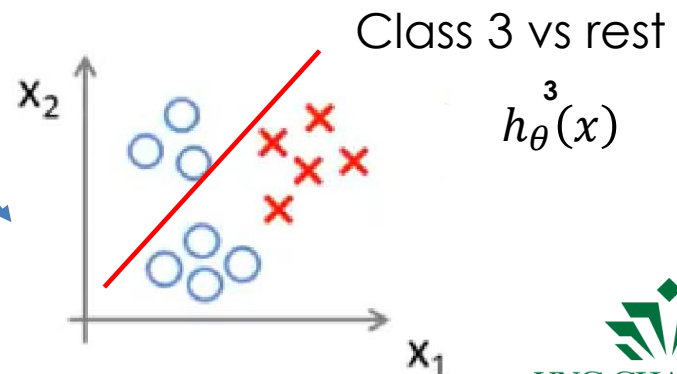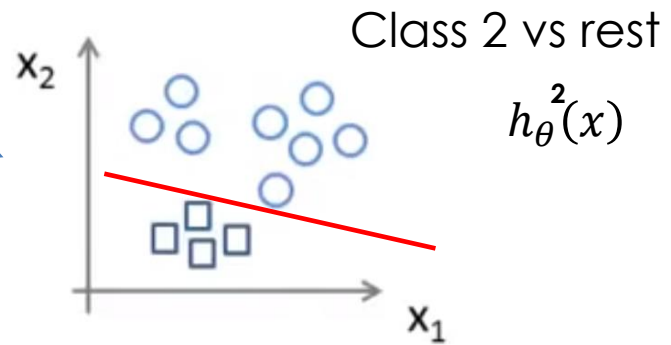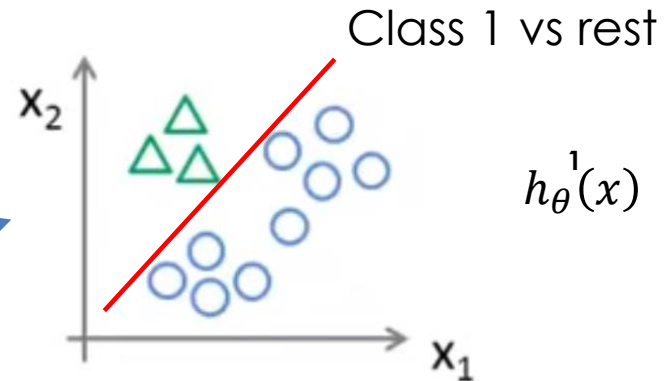
# Logistic regression

**One-vs-all (one-vs-rest):**



Class 1: $\triangle$
Class 2: $\square$
Class 3: $\times$

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

Class 1 vs rest
$h_\theta^1(x)$

Class 2 vs rest
$h_\theta^2(x)$

Class 3 vs rest
$h_\theta^3(x)$

UNC CHARLOTTE

# Logistic regression

**One-vs-all**

Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $y = i$.

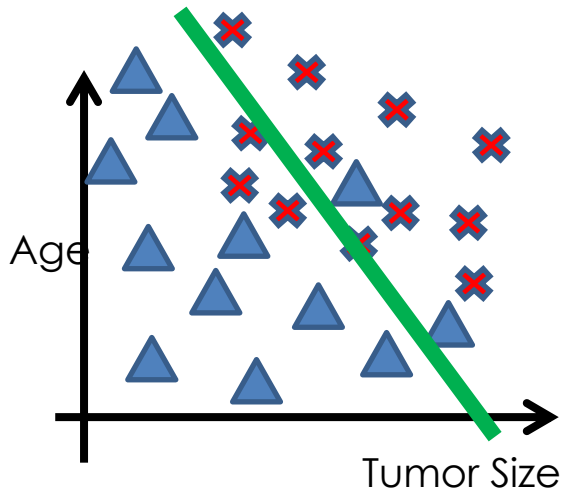On a new input $x$, to make a prediction, pick the class $i$ that maximizes

$$\max_i h_\theta^{(i)}(x)$$

Probability score

UNC CHARLOTTE
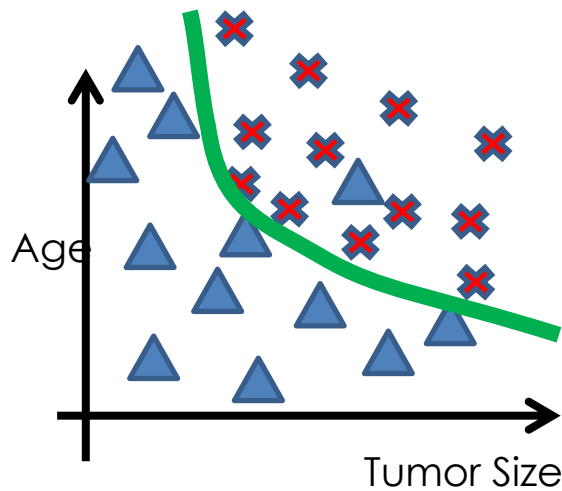
# Logistic regression

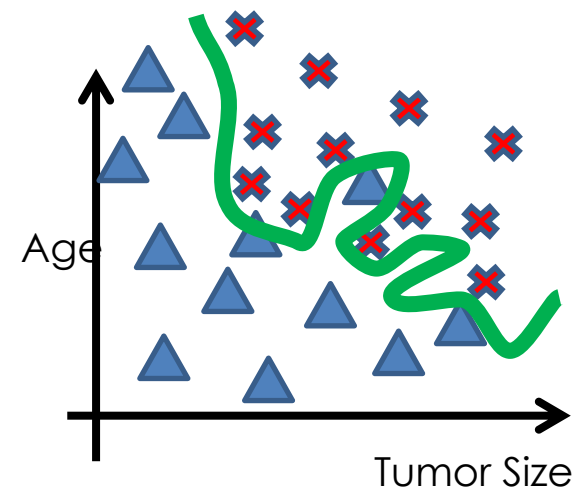- How to perform regularization in logistic regression?

UNC CHARLOTTE

# Overfitting



$$h_\theta(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2)$$

$$h_\theta(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

$$h_\theta(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \theta_6 x_1^3 x_2 + \theta_7 x_1 x_2^3 + \cdots)$$

Underfitting                                                                 Overfitting

◆ Learning the training data too precisely usually leads to poor classification results on new data.

◆ Classifier has to have the ability to generalize.

UNC CHARLOTTE

# Recap: Regularized linear regression

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

$$\min_\theta J(\theta)$$

$n$: Number of features

$\theta_0$ is not panelized

# Regularized logistic regression

- Regularized Logistic Regression

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

$$J_{\text{regularized}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

UNC CHARLOTTE

# Regularized logistic regression

- Regularized Logistic Regression

$$J_{\text{regularized}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

Gradient decent update

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

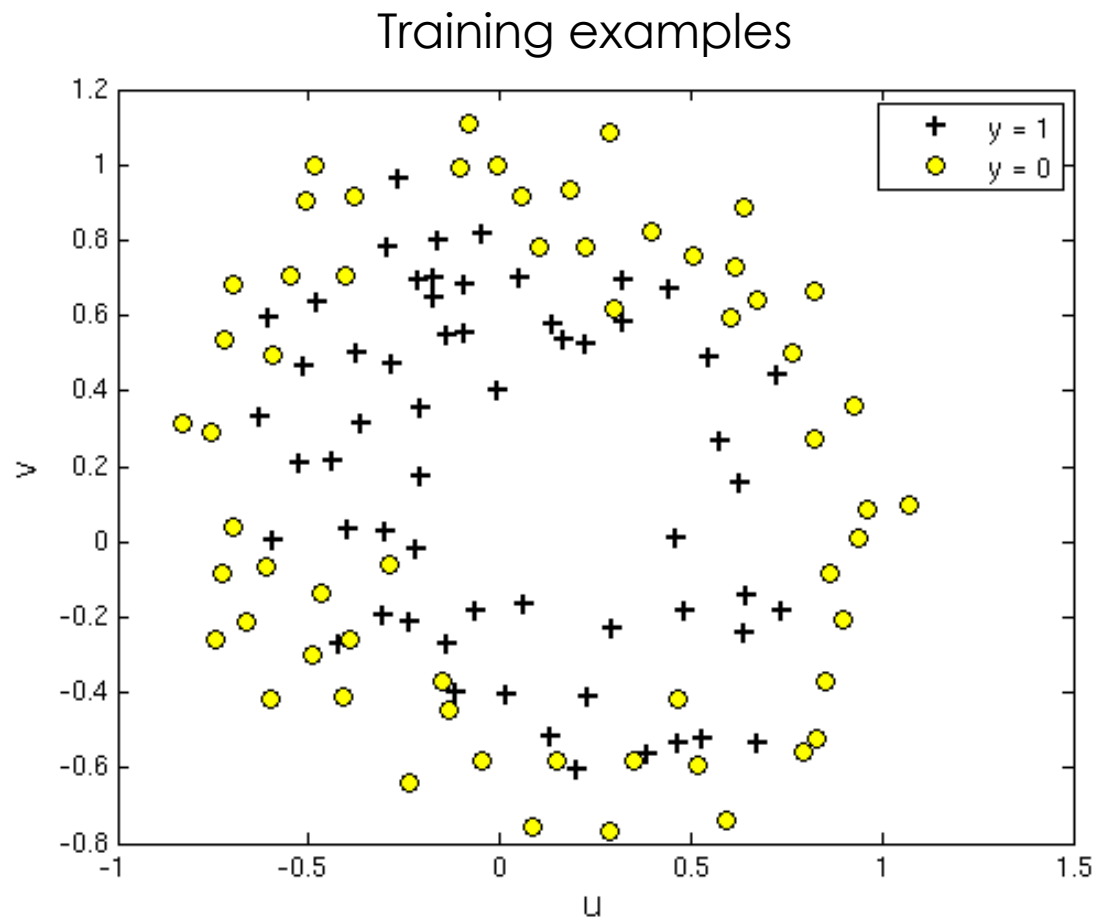<span style="color:red">This looks IDENTICAL to linear regression</span>

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

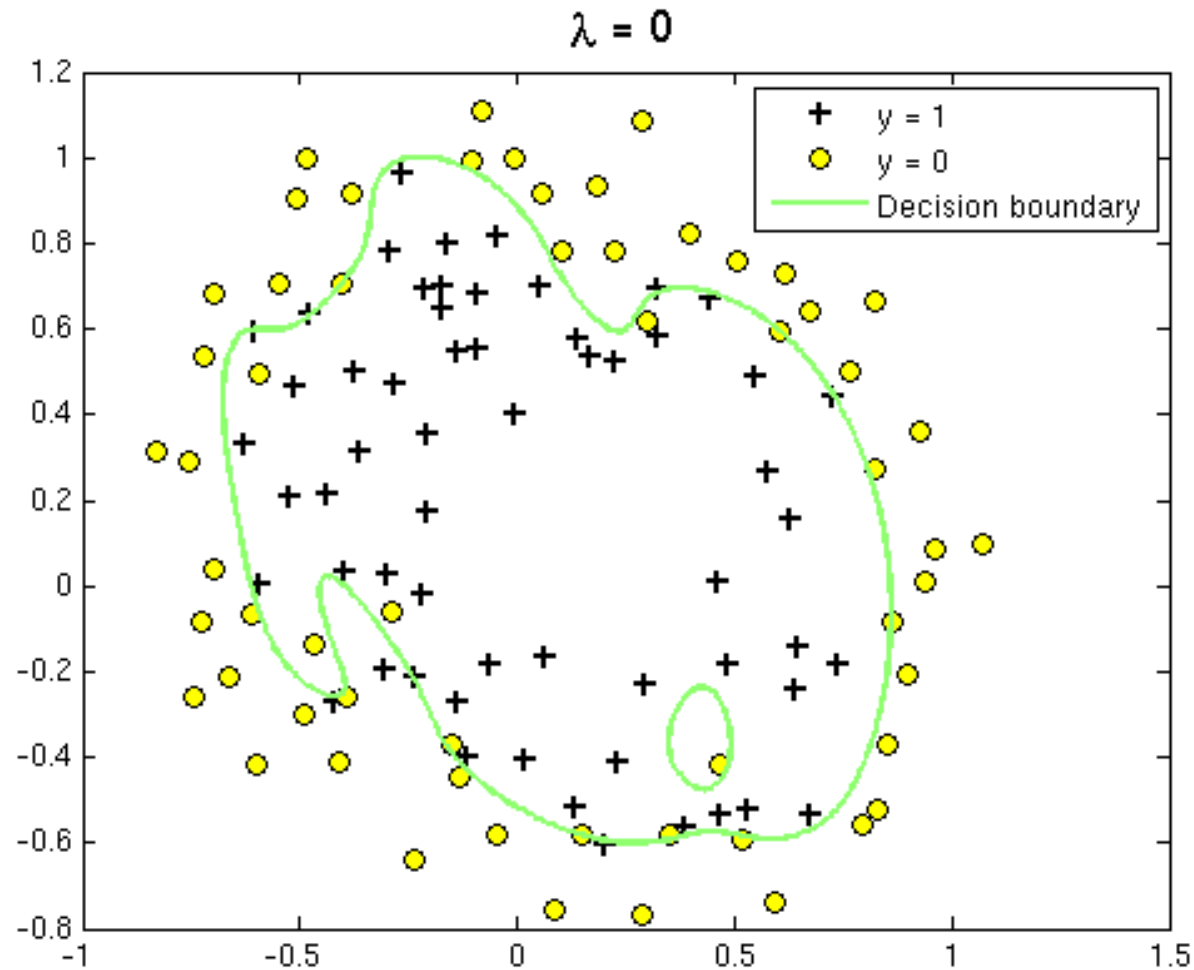However, the form of the model is very different:

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

UNC CHARLOTTE
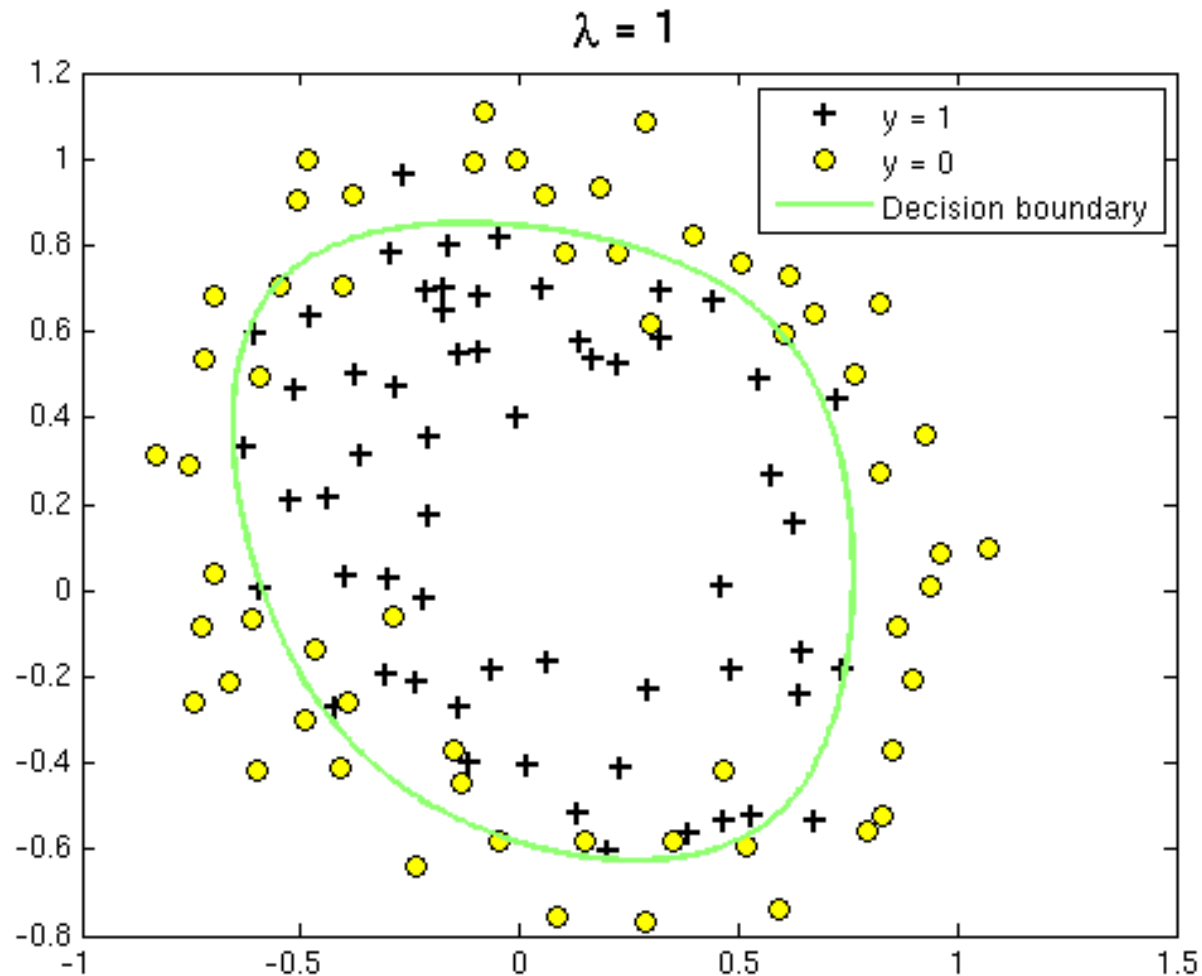
# Regularization

- Example



Training examples

# Regularization

# Regularization

# Regularization



Maybe too much regularization

UNC CHARLOTTE