

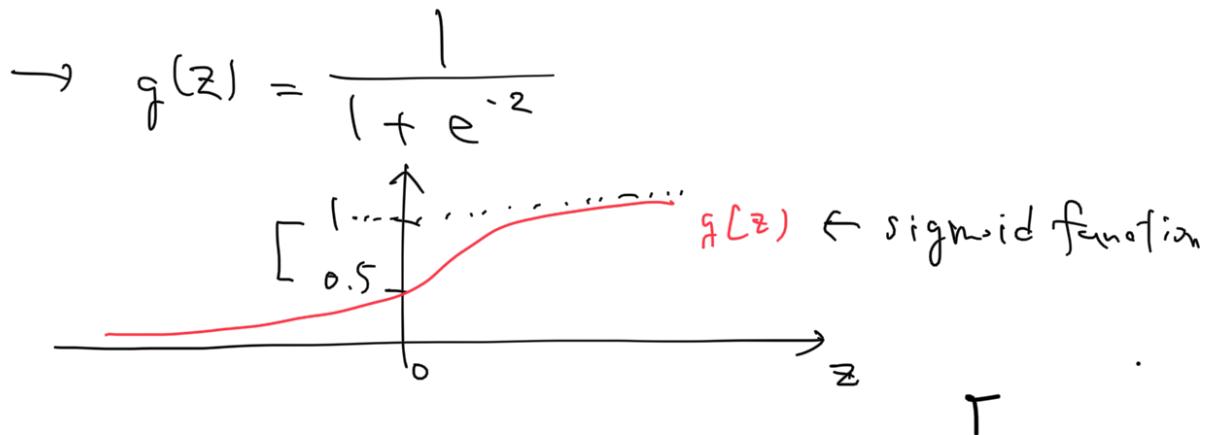


ML w3 logistic regression

|

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

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



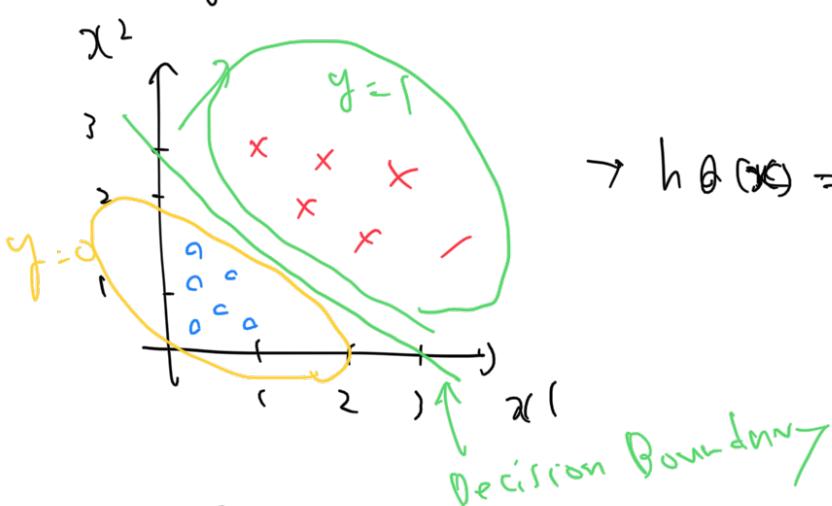
$h_{\theta}(x) = 0.7$... Tell patient that 70% chance of tumor being malignant.

Decision Boundary

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

$$y \geq 1 \quad \text{if } h_{\theta}(x) \geq 0.5 \quad (\textcircled{V})$$

$$y = 0 \quad \text{if } h_{\theta}(x) < 0.5$$



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

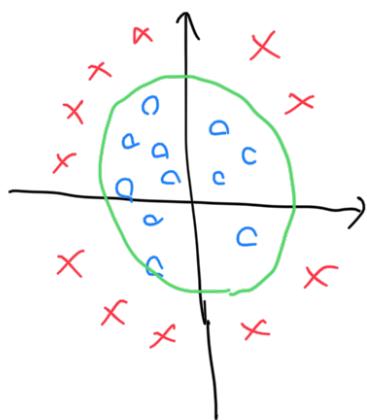
$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\text{Predict } y=1 \text{ if } -3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 = 3$$

(v) $h\theta(x) = 0.5$

Non linear Decision Boundary



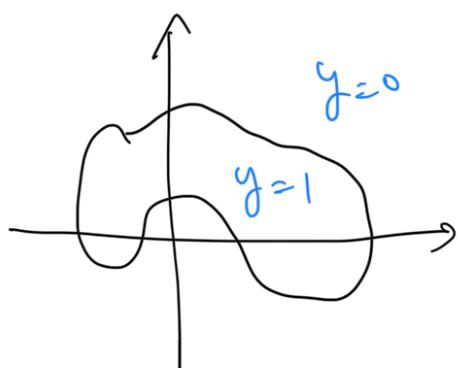
$$h\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Predict $y=1$ if $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 \geq 1$$

$$x_1^2 + x_2^2 = 1$$



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

Cost Function

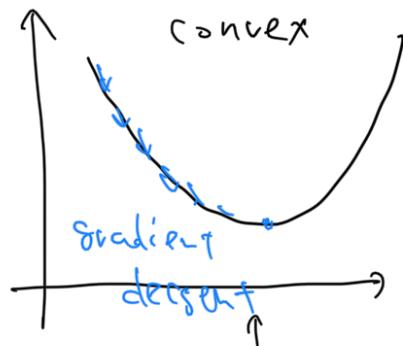
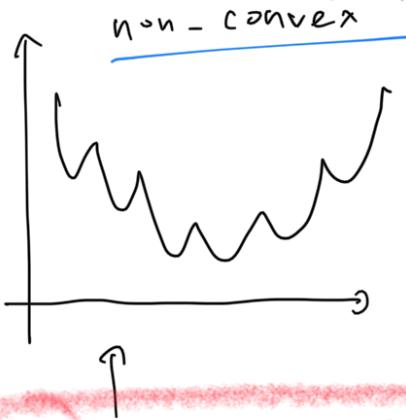
- Training Set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$
- m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}_{\mathbb{R}^{n+1}}$ $x_0 = 1, y \in \{0, 1\}$

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

(Linear Regression model)

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

cost($h_{\theta}(x^{(i)})$, $y^{(i)}$) ← squared error



(Logistic Regression)

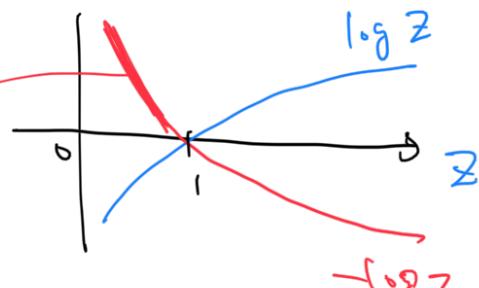
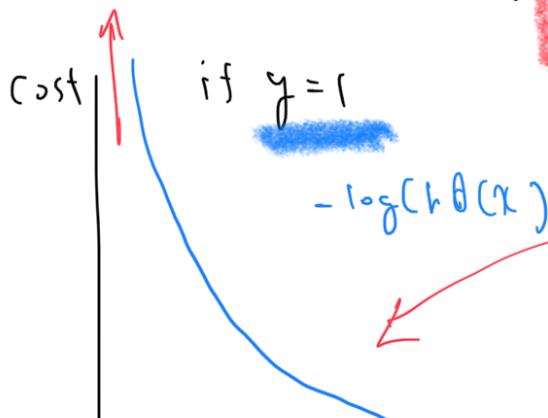
$$\text{Cost}(h_{\theta}(x), y^{(i)}) = \frac{1}{1 + e^{-\theta^T x}}$$

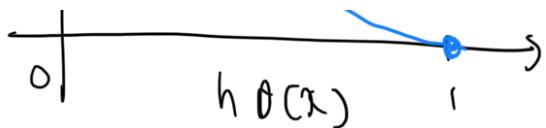
non-convex function
非凸関数

非線形、
局所最適。

Convex, gradient descent → 局所最適

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



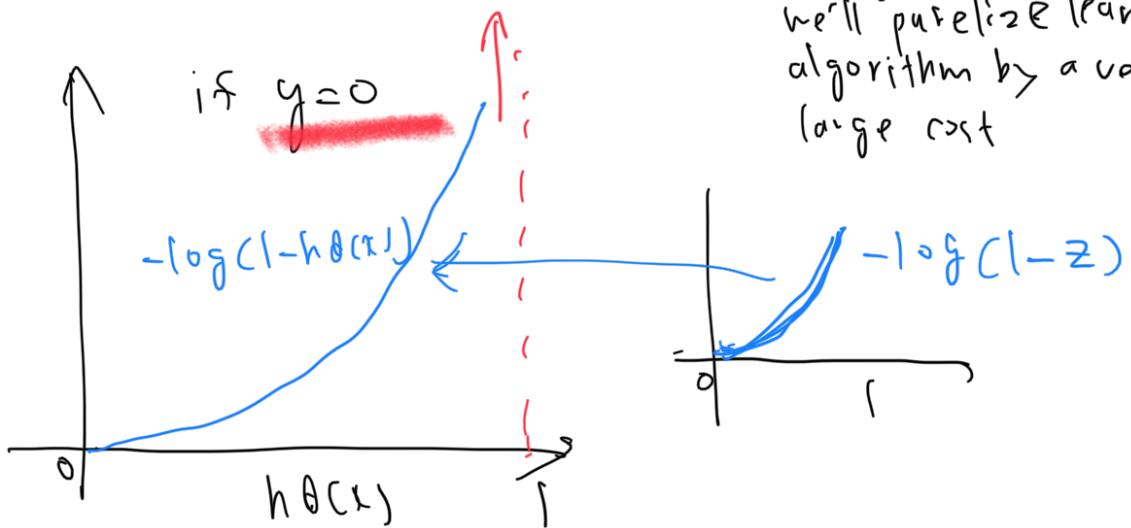


Cost = 0 if $y=1$, $h_\theta(x)=1$

But as $h_\theta(y) \rightarrow 0$
cost $\rightarrow \infty$

if $h_\theta \rightarrow 0$,
(predict $P(y=1|x;\theta)=0$)

but $y=1$,
we'll penalize learning
algorithm by a very
large cost



Simplified Cost Function and Gradient Descent

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

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

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

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

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

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

Minimize

$$\min \theta J(\theta)$$

$$h_{\theta}(x) = \theta^T x \in \text{Linear Regression}$$

Repeat {

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

}

simultaneously update

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

↳ Logistic Regression

Vectorized Implementation

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$$

Advanced Optimization

Gradient Descent General Form

Repeat {

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

}

[Conjugate/BFGS ...] \hookrightarrow

pros

No need to manually pick α

Example

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$\rightarrow J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

```
function [jVal, gradient]
    = costFunction(theta)
        jVal = (theta(1) - 5)^2 ..
```

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

$$(\text{theta}(1) - 5) \\ \text{gradient} = 2 * \text{res}(2, 1) \\ \text{gradient}(1) = 2 * (\text{theta}(0) - 5) \\ \text{gradient}(2) = 2 * (\text{theta}(1) - 5)$$

... advanced: ex: minimize (@costFunc(theta, initialTheta, options))

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \begin{array}{l} \text{theta}(1) \\ \text{theta}(1) \\ \vdots \\ \text{theta}(n) \end{array}$$

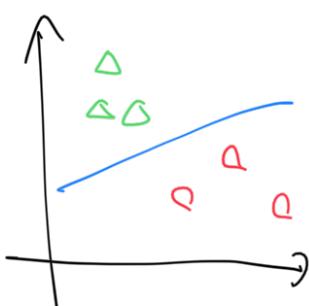
function [jVal, gradient] = costFunction(theta).

$$\text{gradient}(1) \\ = (1) \\ = \dots (n+1)$$

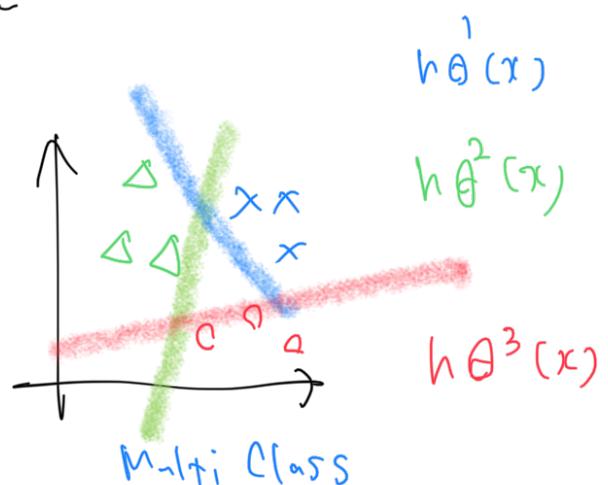
Multiclass Classification One-vs-all

e.g. Sunny, Cloudy, Rainy ..

$y = 1 \quad 2 \quad 3$



Binary



Multi Class

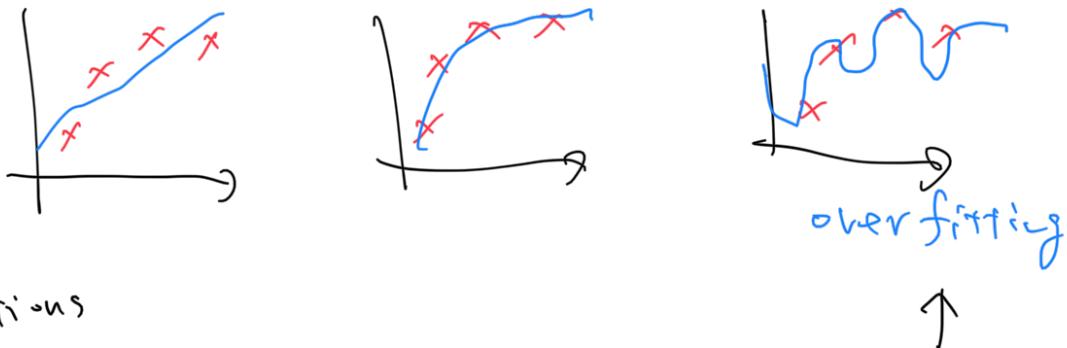
$$h_{\theta}^{(i)}(x) = P(y=i|x; \theta)$$

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

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

Regularization — solving the problem of overfitting.



Options

1. Reduce number of features

2. Regularization

Cost Function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \underbrace{\sum_{j=1}^n \theta_j^2}_{\text{penalize}}$$

$$\theta_n \rightarrow \theta_0$$

because x^n and λ
is very small.

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

← set
not too
large value.

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Regularized Linear Regression

Gradient descent

$$\begin{aligned}
 & \text{Repeat} \quad \left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \end{array} \right. \\
 & \theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \left. \begin{array}{l} \text{if } \theta_j \neq 0 \\ \text{else } 0.99 \dots \end{array} \right. \\
 & 1 - \alpha \frac{\lambda}{m} < 1
 \end{aligned}$$

Normal Equation

$$X = \begin{bmatrix} x^{(1)\top} \\ \vdots \\ \vdots \\ x^{(m)\top} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

$$\rightarrow \theta = (X^\top X + \lambda \begin{bmatrix} 0 & & \\ & \ddots & \\ & & 1 \end{bmatrix})^{-1} X^\top y$$

$n+1$ by $n+1$
matrix

Non-invertibility (optional / advanced)

Suppose ($m < n$)

$$\theta = (X^\top X)^{-1} X^\top y \quad \text{pinv}$$

if $\lambda > 0$, $X^\top X$ is $n \times n$ and invertible

$$\theta = \underbrace{(X^T X + \lambda \begin{bmatrix} 0 & & \\ & \ddots & \\ & & 1 \end{bmatrix})^{-1}}_{\uparrow \text{non invertible}} X^T y$$

Regularized logistic Regression

- Cost function (see above)
- Gradient Descent
- Advanced Optimization.