

Math 7243

Machine Learning and Statistical Learning Theory

Section 3. Logistic Regression

Instructor: He Wang

Department of Mathematics

Northeastern University

Review :

1. Linear regression
2. Locally weighted linear regression
3. Ridge Regression
4. Lasso Regression

Training Data:

$$D = \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(n)}, y^{(n)})\}$$

features *labels*

Assumption: Linear Regression

$$h(\vec{x}) = \vec{\theta}^T \vec{x} = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

Today:

1. Logistic Regression (binary)
2. Softmax Regression (multiclass)

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

- Example: Data of students sleep time, study time, and pass/fail.

If we know the test scores, we can use linear regression to predict the test scores.

$$P(Y=1 | x_1=7, x_2=7) = ?$$

$$P(Y=0 | x_1=2, x_2=2) = ?$$

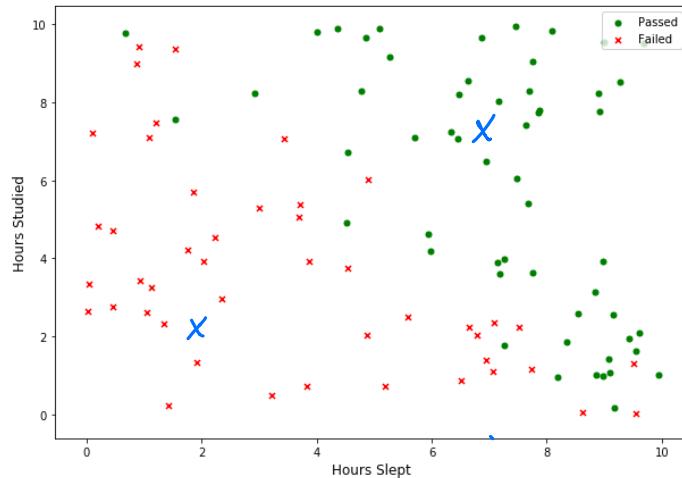
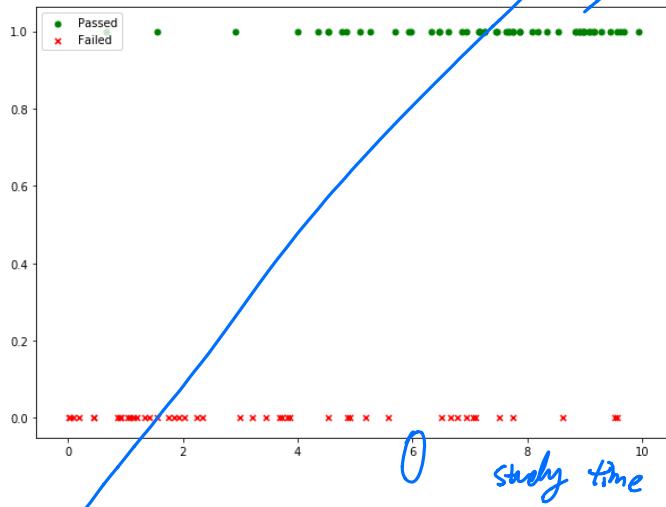
➤ Logistic regression

$$h(\vec{x}) = P(Y=1 | \vec{x})$$

Logistic regression is a classification algorithm, used to predict probabilities based on given set of independent variables.

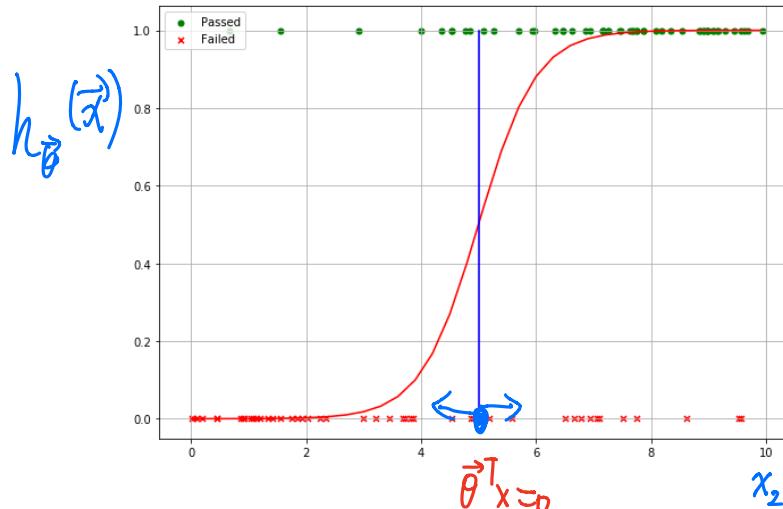
$$P(Y=1 | x_2=6) = ?$$

(linear) regression



x_1	x_2	Y
7.63	7.40	1
2.03	3.93	0
3.82	0.72	0
7.15	3.89	1
6.47	8.19	1
...

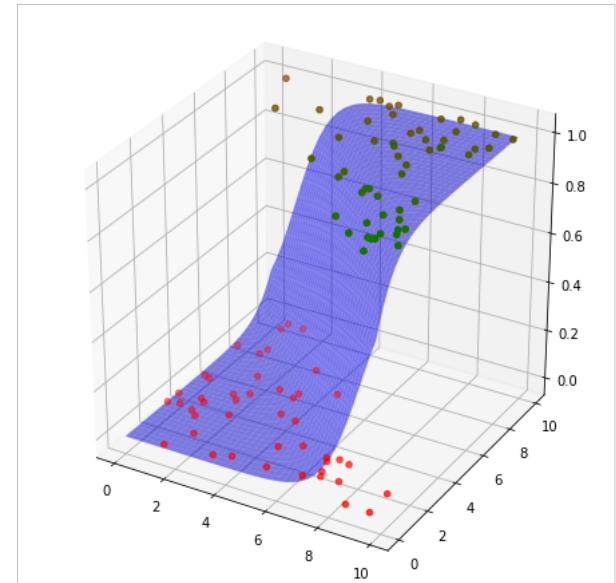
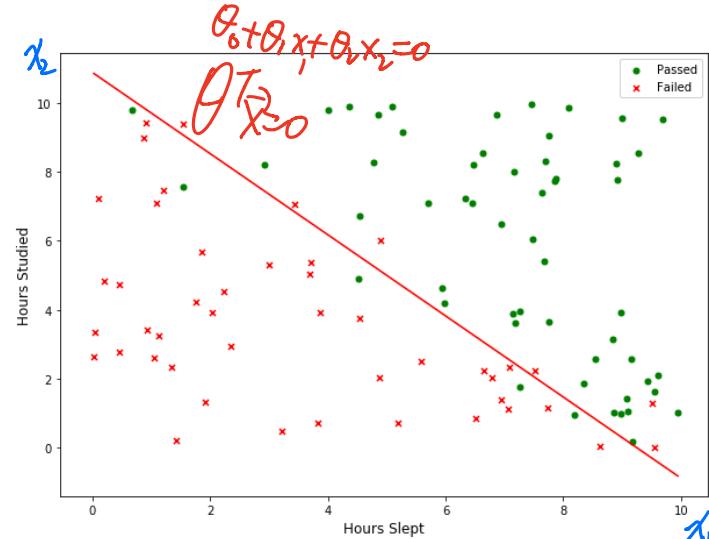
Scenes
Z
90
50
40
80
90
;



$$h(\vec{x}) = \text{at } \vec{x} \in \{ \vec{x} \in \mathbb{R}^d \mid f(\vec{x}) = 0 \}$$

➤ Decision Boundary $\vec{\theta}^T \vec{x} = f(\vec{x}) = 0$

logistic regression prediction function returns a probability between 0 and 1, in order to predict which class this data belongs we need to set a threshold.



➤ Sigmoid function.

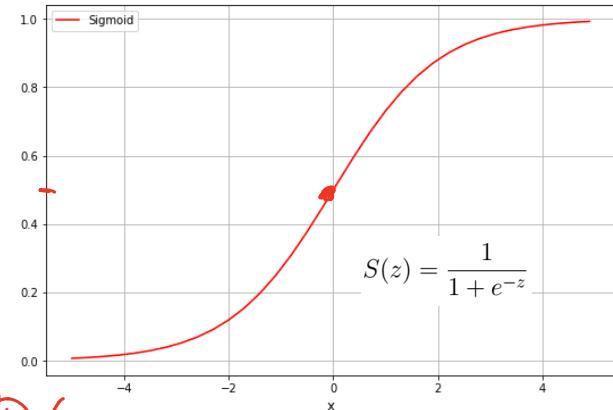
$$S(z) = \frac{1}{1 + e^{-z}}$$

$$S(0)=0.5$$

The sigmoid function maps any real value into a value in $[0,1]$.

Hypothesis

Goal: $h_{\vec{\theta}}(\vec{x}) = S(\vec{\theta}^T \vec{x}) = \frac{1}{1 + e^{-\vec{\theta}^T \vec{x}}} = \text{probability}$



$$P(y=1 | \vec{x})$$

class Prediction:

$$y(\vec{x}) = \begin{cases} \text{1-pass} & \text{if } h(\vec{x}) \geq 0.5 \\ \text{0-fail} & \text{if } h(\vec{x}) < 0.5 \end{cases}$$

Boundary $\vec{\theta}^T \vec{x} = 0$

Remark: $\log \left(\frac{h(\vec{x})}{1-h(\vec{x})} \right) = \vec{\theta}^T \vec{x}$

$$= \log \left(\frac{P(y=1 | \vec{x})}{P(y=0 | \vec{x})} \right)$$

• other functions

① Step function $s(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$

② Sign function $s(z) = \begin{cases} -1 & z < 0 \\ 1 & z > 0 \end{cases}$

③ Tanh $s(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

Assumption: $S(\theta^T \vec{x}) = \frac{1}{1 + e^{-\theta^T \vec{x}}}$

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

$$\begin{cases} P(y=1|\vec{x};\theta) = h_\theta(\vec{x}) \\ P(y=0|\vec{x};\theta) = 1 - h_\theta(\vec{x}) \end{cases} \Rightarrow P(y|\vec{x};\theta) = h_\theta(\vec{x})^y (1 - h_\theta(\vec{x}))^{1-y}$$

Maximize Likelihood function:

$$L(\theta) = P(\vec{y} | X; \theta)$$

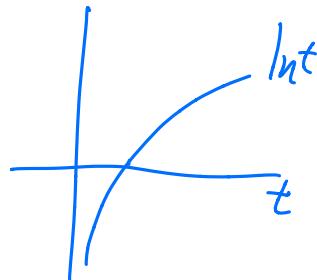
$$= \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) \quad (\text{i.i.d.})$$

$$= \prod_{i=1}^n h(x^{(i)})^{y^{(i)}} (1 - h(x^{(i)}))^{1-y^{(i)}}$$

Maximize log likelihood:

$$l(\theta) = \ln(L(\theta))$$

$$= \sum_{i=1}^n \left(\underbrace{y^{(i)} \ln h(x^{(i)})}_{<0} + \underbrace{(1-y^{(i)}) \ln (1-h(x^{(i)}))}_{<0} \right)$$



log cost for one observation:

$$\text{cost}(\theta) = \begin{cases} -\ln(h(x^{(i)})) & \text{if } y^{(i)} = 1 \\ -\ln(1 - h(x^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$= -(y^{(i)} \ln(h(x^{(i)})) + (1-y^{(i)}) \ln(1-h(x^{(i)})))$$

log cost function:

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} \ln h(x^{(i)}) + (1-y^{(i)}) \ln (1-h(x^{(i)})))$$

- Gradient: $\nabla J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_d} \end{bmatrix}$

$$h(x^{(i)}) = S(\theta^T x^{(i)})$$

$$S'(z) = S(z) (1 - S(z))$$

$$\frac{\partial h(x^{(i)})}{\partial \theta_j} = S(z) (1 - S(z)) \cdot x_j^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{n} \sum_{i=1}^n \left(y^{(i)} \frac{1}{S(z)} \cancel{S(z)} (1 - S(z)) \cdot x_j^{(i)} - (1-y^{(i)}) \frac{1}{1-S(z)} \cancel{S(z)} (1 - S(z)) \cdot x_j^{(i)} \right)$$

$$= \frac{1}{n} \sum_{i=1}^n (S(\theta^T x^{(i)}) - y^{(i)}) x_j^{(i)} = \frac{1}{n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

- Gradient descent: $\vec{\theta}^{\text{next}} = \vec{\theta} - \alpha \nabla J(\vec{\theta})$

- Newton's method: $\vec{\theta}^{\text{next}} = \vec{\theta} - H^{-1} \nabla J(\vec{\theta})$

Hessian $H_{jk} = \frac{\partial^2 J}{\partial \theta_j \partial \theta_k} = \frac{1}{n} \sum_{i=1}^n h(x^{(i)}) (1 - h(x^{(i)})) x_j^{(i)} x_k^{(i)}$

Question: If $y \in \{-1, 1\}$, assume
 $\text{find } \nabla J(\theta).$

- $P(y|\vec{x}; \theta) = h_{\theta}(\vec{x})^{\frac{y}{2}} (1 - h_{\theta}(\vec{x}))^{\frac{1-y}{2}}$

better $\curvearrowleft h_{\theta}(y\vec{x})$

\downarrow

- $J(\theta) = -\frac{1}{n} \sum_{i=1}^n \ln(h(y\vec{x}))$

$$= \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^{(i)} \theta^T x^{(i)}})$$

- Gradient ∇J :

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

$$\begin{cases} P(y=1 | \vec{x}, \theta) = h_{\theta}(\vec{x}) \\ P(y=-1 | \vec{x}, \theta) = 1 - h_{\theta}(\vec{x}) \end{cases}$$

To use Newton's method $\theta^{\text{next}} = \theta - H^{-1} \nabla J$

- Hessian matrix H

$$H_{jk} = \frac{\partial^2 J}{\partial \theta_j \partial \theta_k}$$

$$= \frac{1}{n} \sum_{i=1}^n h(y^{(i)} x^{(i)}) (1 - h(y^{(i)} x^{(i)})) y^{(i)} x_j^{(i)} y^{(i)} x_k^{(i)}$$

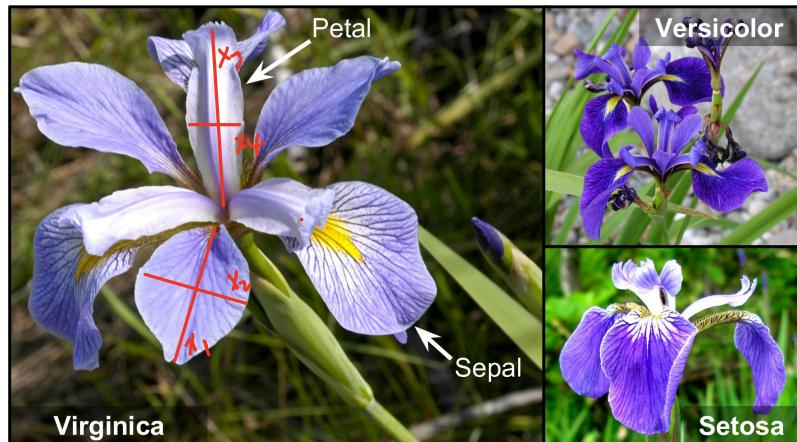
$$= \frac{1}{n} \sum_{i=1}^n h(y^{(i)} x^{(i)}) (1 - h(y^{(i)} x^{(i)})) x_j^{(i)} x_k^{(i)}$$

Gradient descent: $\theta^{\text{next}} = \theta - \alpha \nabla J(\theta)$

➤ Softmax Regression (Multinomial Logistic Regression)

- Flowers of three iris plant species:

The famous Iris database, first used by Sir R.A. Fisher(1936), is best known database to be found in the pattern recognition literature. It contains the **sepal** and **petal length** and **width** of 150 iris flowers of three different species: Iris-Setosa, Iris-Versicolor, and Iris-Virginica.



$$\vec{x} \in \mathbb{R}^4$$

Data features: sepal length, sepal width, petal length, petal width.

$$y \in \mathbb{Z}_3^3 = \{0, 1, 2\}$$

Classes: 0-Iris-Setosa, 1-Iris-Versicolour, 2-Iris-Virginica

$$y \in \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$$

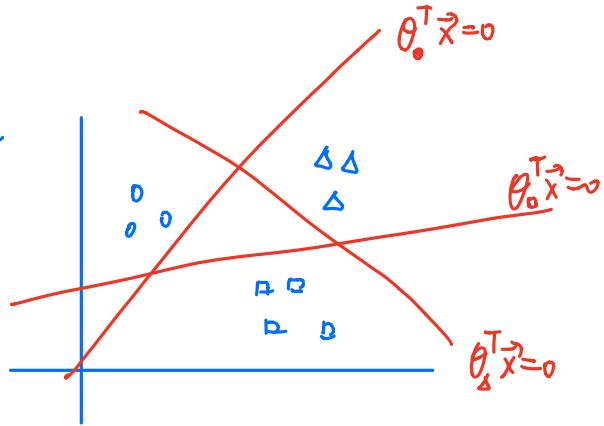
or $\left\{ \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$

Want $\vec{\theta}_0, \vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^4$

$$\vec{\theta} = \begin{bmatrix} \vec{\theta}_0 & \vec{\theta}_1 & \vec{\theta}_2 \end{bmatrix}_{4 \times 3}$$

x_1	x_2	x_3	x_4	y
[5.1, 3.5, 1.4, 0.2],				0
[4.9, 3. , 1.4, 0.2],				0
[4.7, 3.2, 1.3, 0.2],				0
[4.6, 3.1, 1.5, 0.2],				1
[5. , 3.6, 1.4, 0.2],				1
[5.4, 3.9, 1.7, 0.4],				2
[4.6, 3.4, 1.4, 0.3],				1
[5. , 3.4, 1.5, 0.2],				1
[4.4, 2.9, 1.4, 0.2]				2
...				1

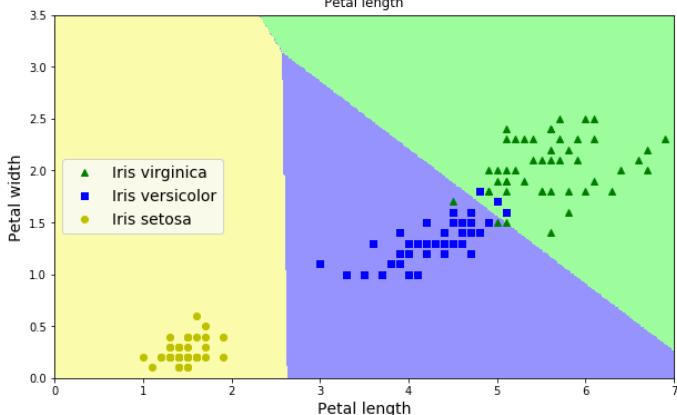
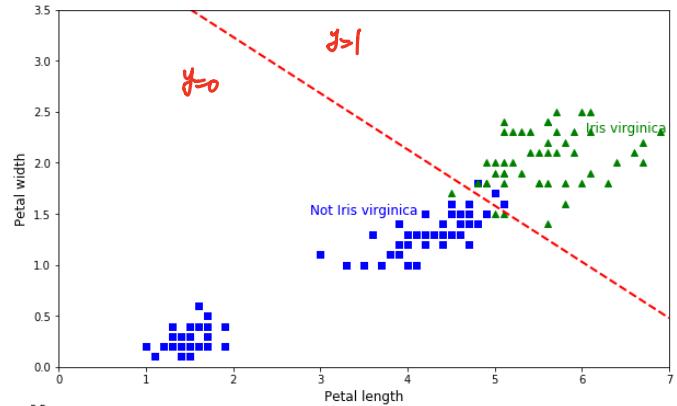
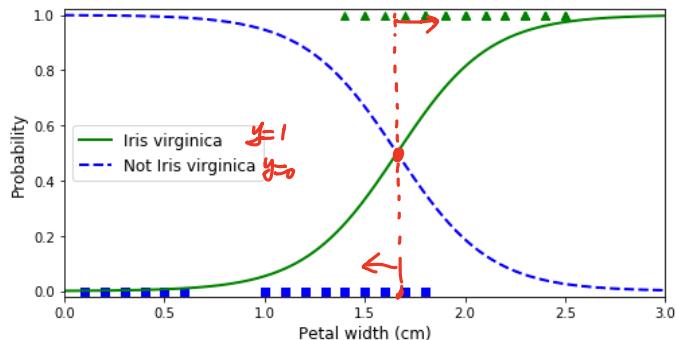
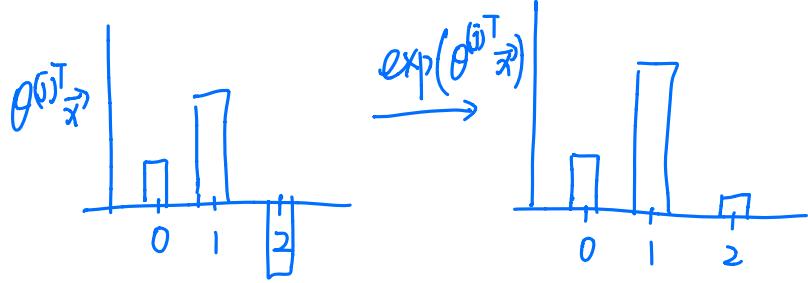
such that



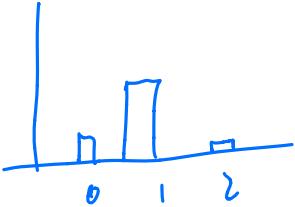
Hypothesis:

$$h_{\theta}(\vec{x}) = \begin{bmatrix} P(y=0|\vec{x}; \theta) \\ P(y=1|\vec{x}; \theta) \\ P(y=2|\vec{x}; \theta) \end{bmatrix} = \frac{1}{\sum_{j=0}^2 \exp(\theta^{(j)\top} \vec{x})} \begin{bmatrix} \exp(\theta^{(0)\top} \vec{x}) \\ \exp(\theta^{(1)\top} \vec{x}) \\ \exp(\theta^{(2)\top} \vec{x}) \end{bmatrix}$$

↑
distribution vector



$$\rightarrow \frac{\exp(\theta^{(i)\top} \vec{x})}{\sum_{j=0}^2 \exp(\theta^{(j)\top} \vec{x})}$$



log. Cost Function: (Cross entropy cost function)

$$H(p, q) = - \sum p(x) \log(q(x))$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left(1\{y^{(i)}=0\} \ln(P(y^{(i)}=0 | \vec{x}, \theta)) + 1\{y^{(i)}=1\} \ln(P(y^{(i)}=1 | \vec{x}^{(i)}, \theta)) + 1\{y^{(i)}=2\} \ln(P(y^{(i)}=2 | \vec{x}^{(i)}, \theta)) \right)$$

$$= -\frac{1}{n} \sum_{i=1}^n \sum_{k=0}^2 \left(1\{y^{(i)}=k\} \ln \left(\frac{\exp(\theta^{(k)\top} \vec{x}^{(i)})}{\sum_{j=0}^2 \exp(\theta^{(j)\top} \vec{x}^{(i)})} \right) \right)$$

$h(\vec{x}^{(i)})^{(k)}$

For the case of binary classification

$$h_\theta(\vec{x}) = \frac{1}{\exp(\theta^{(0)\top} \vec{x}) + \exp(\theta^{(1)\top} \vec{x})} \begin{bmatrix} \exp(\theta^{(0)\top} \vec{x}) \\ \exp(\theta^{(1)\top} \vec{x}) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{1 + \exp((\theta^{(0)\top} - \theta^{(1)\top}) \vec{x})} \\ 1 - \frac{1}{1 + \exp((\theta^{(0)\top} - \theta^{(1)\top}) \vec{x})} \end{bmatrix}$$

Set $\theta = \theta^{(0)} - \theta^{(1)}$

$$= \begin{bmatrix} \frac{1}{1 + \exp(-\theta^\top \vec{x})} \\ 1 - \frac{1}{1 + \exp(-\theta^\top \vec{x})} \end{bmatrix}$$

where $1\{ \cdot \}$ is the "indicator function"

$$\begin{cases} 1\{\text{true}\} = 1 \\ 1\{\text{false}\} = 0 \end{cases}$$

- Use gradient descent or Newton's method to find $\theta^{(0)}, \theta^{(1)}, \theta^{(2)} \dots$
? H is singular.
- Gradient $\nabla_{\theta^{(k)}} J(\theta) = \frac{1}{n} \sum_{i=1}^n (h(x^{(i)})^{(k)} - 1(y^{(i)}=k)) x^{(i)}$