

Lecture 3

Logistic Regression & Softmax Regression

Rui Xia

School of Computer Science & Engineering
Nanjing University of Science & Technology

<http://www.nustm.cn/~rxia>

Supervised Learning

- Regression



Share Price
"\$ 24.50"

Continuous Labels
Regression

- Classification

Feature Space \mathcal{X}

Words in a document

Label Space \mathcal{Y} :

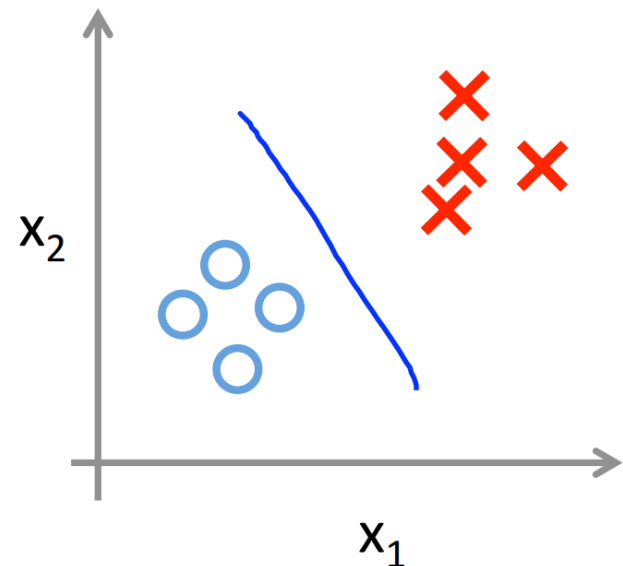
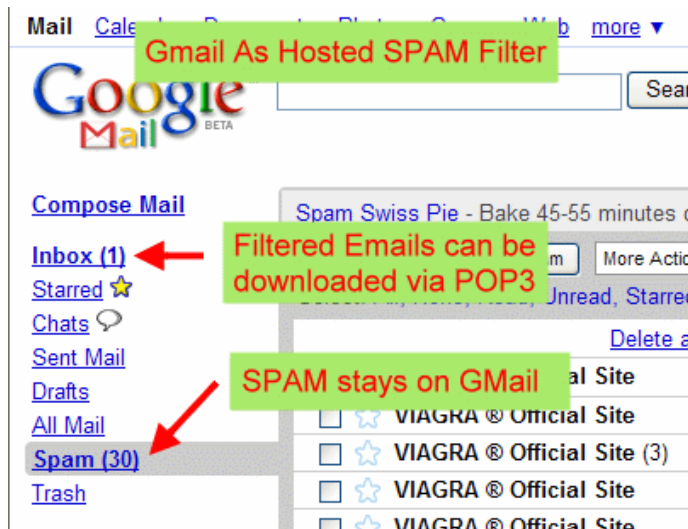
"Sports"
"News"
"Science"
...

Discrete Labels
Classification

Logistic Regression

Introduction

- Logistic Regression is a **classification** model, although it is called “regression”;
- Logistic regression is a binary classification model;
- Logistic regression is a linear classification model. It has a linear decision boundary (hyperplane), but with a nonlinear activation function (Sigmoid function) to model the posterior probability.

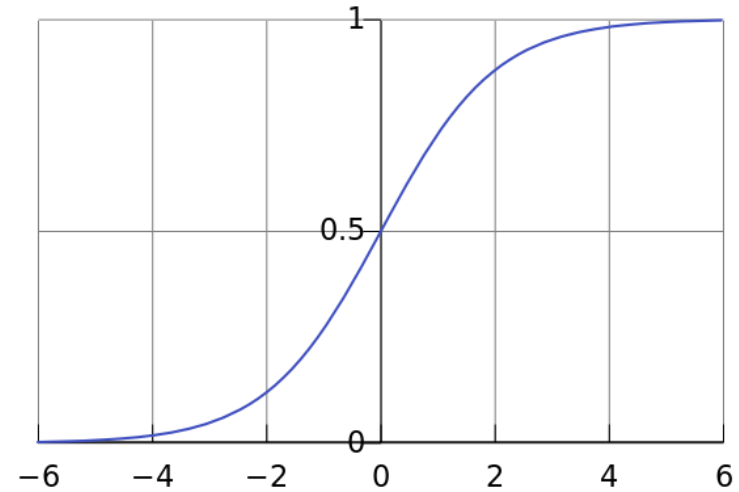


Model Hypothesis

- Sigmoid Function

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

$$\frac{d\delta(z)}{dz} = \delta(z) (1 - \delta(z))$$



- Hypothesis

$$p(y = 1 | \mathbf{x}; \boldsymbol{\theta}) = h_{\boldsymbol{\theta}}(\mathbf{x}) = \delta(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

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

- Hypothesis (Compact Form)

$$p(y | \mathbf{x}; \boldsymbol{\theta}) = (h_{\boldsymbol{\theta}}(\mathbf{x}))^y (1 - h_{\boldsymbol{\theta}}(\mathbf{x}))^{(1-y)} = \left(\frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}} \right)^y \left(1 - \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}} \right)^{(1-y)}$$

Learning Algorithm

- (Conditional) Likelihood Function

$$\begin{aligned} L(\boldsymbol{\theta}) &= \prod_{i=1}^N p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \prod_{i=1}^N \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^{y^{(i)}} \left(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^{(1-y^{(i)})} \\ &= \prod_{i=1}^N \left(\frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}} \right)^{y^{(i)}} \left(1 - \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}^{(i)}}} \right)^{(1-y^{(i)})} \end{aligned}$$

- Maximum Likelihood Estimation

$$\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \Leftrightarrow \max_{\boldsymbol{\theta}} \sum_{i=1}^N y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))$$

The neg log-likelihood function is also known as the **Cross-Entropy** cost function

Unconstraint Optimization

- Unconstraint Optimization Problem

$$\max_{\theta} \sum_{i=1}^N y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(\mathbf{x}^{(i)}))$$

- Optimization Methods
 - Gradient Descent
 - Stochastic Gradient Descent
 - Newton Method
 - Quasi-Newton Method
 - Conjugate Gradient
 - ...

Gradient Descent/Ascent

- Gradient Computation

$$\begin{aligned}\frac{dl(\boldsymbol{\theta})}{d\boldsymbol{\theta}} &= \sum_{i=1}^N \left(y^{(i)} \frac{1}{h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} \right) \frac{d}{d\boldsymbol{\theta}} h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \\ &= \sum_{i=1}^N \left(y^{(i)} \frac{1}{h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} \right) h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \frac{d}{d\boldsymbol{\theta}} \boldsymbol{\theta}^T \mathbf{x}^{(i)} \\ &= \sum_{i=1}^N \left(y^{(i)} (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) - (1 - y^{(i)}) h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)} \\ &= \sum_{i=1}^N \boxed{(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \mathbf{x}^{(i)}} \quad \text{Error} \times \text{Feature}\end{aligned}$$

- Gradient Ascent Optimization

$$\boldsymbol{\theta} := \boldsymbol{\theta} + \alpha \sum_{i=1}^N (y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \mathbf{x}^{(i)}$$

Stochastic Gradient Descent

- Randomly choose a training sample

$$(x, y)$$

- Compute gradient

$$(y - h_{\theta}(x))x$$

- Updating weights

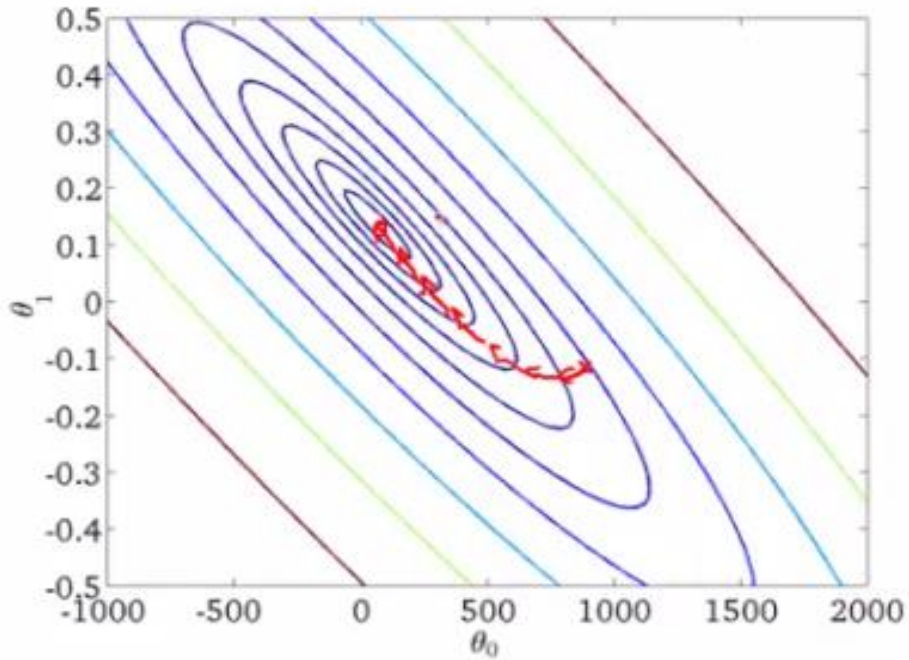
$$\theta := \theta + \alpha(y - h_{\theta}(x))x$$

- Repeat...

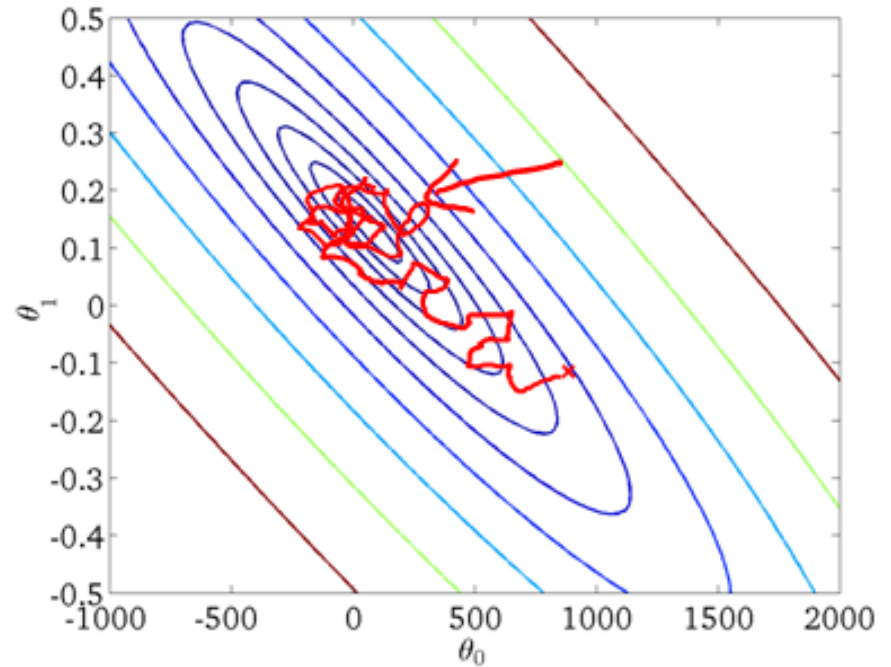
Gradient descent -- **batch** updating

Stochastic gradient descent -- **online** updating

GD vs. SGD



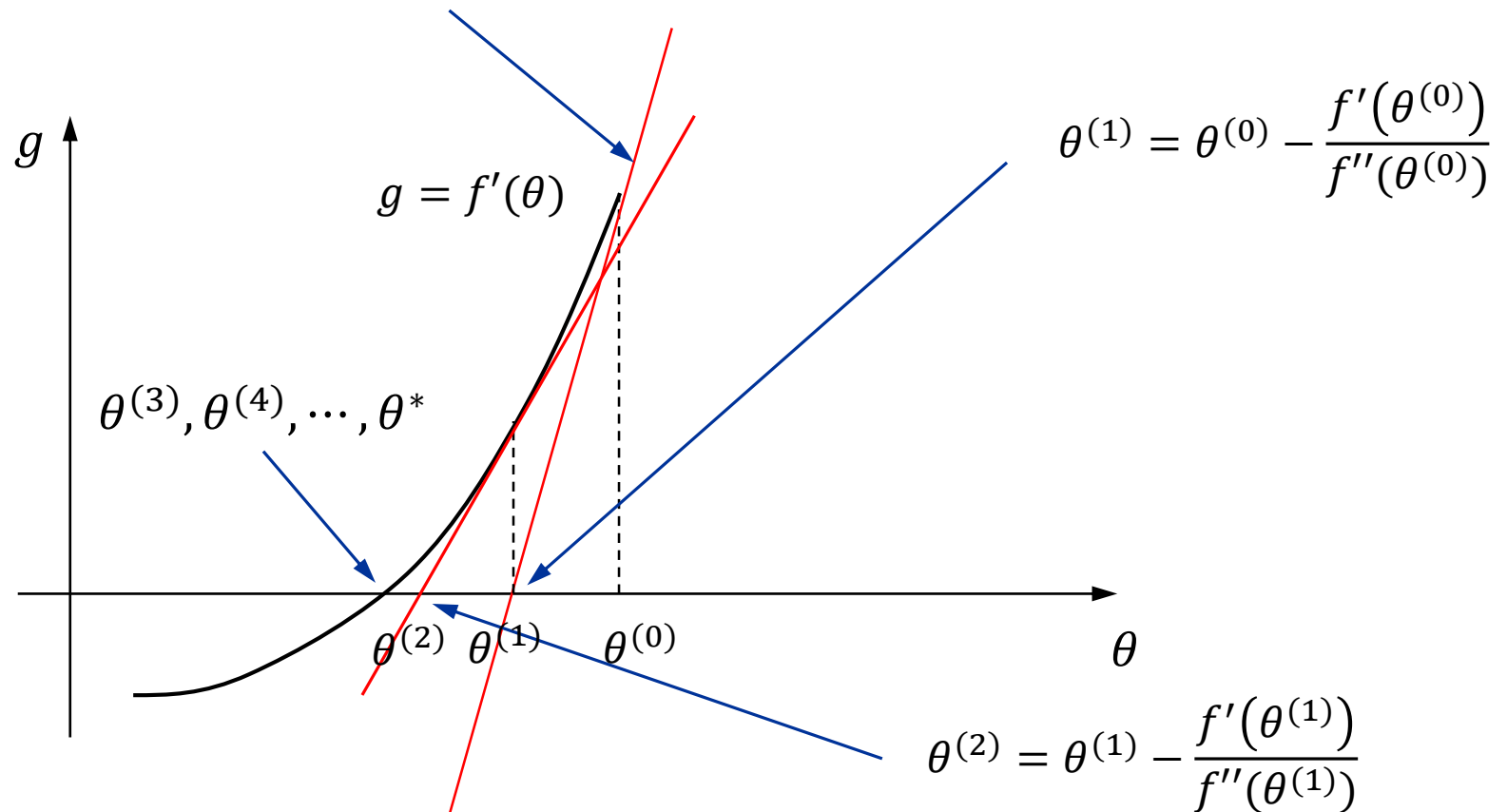
Gradient Descent (GD)



Stochastic Gradient Descent (SGD)

Illustration of Newton's Method

tangent line: $g = f'(\theta^{(0)}) + f''(\theta^{(0)})(\theta - \theta^{(0)})$



Newton's Method

- Problem

$$\arg \min f(\boldsymbol{\theta}) \Leftrightarrow \text{solve} : \nabla f(\boldsymbol{\theta}) = 0$$

- Second-order Taylor expansion

$$\phi(\boldsymbol{\theta}) = f(\boldsymbol{\theta}^{(k)}) + \nabla f(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}) + \frac{1}{2} \nabla^2 f(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)})^2 \approx f(\boldsymbol{\theta})$$

$$\nabla \phi(\boldsymbol{\theta}) = 0 \Rightarrow \boldsymbol{\theta} = \boldsymbol{\theta}^{(k)} - \nabla^2 f(\boldsymbol{\theta}^{(k)})^{-1} \nabla f(\boldsymbol{\theta}^{(k)})$$

- Newton's method (also called Newton-Raphson method)

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \boxed{\nabla^2 f(\boldsymbol{\theta}^{(k)})}^{-1} \nabla f(\boldsymbol{\theta}^{(k)})$$

Hessian Matrix

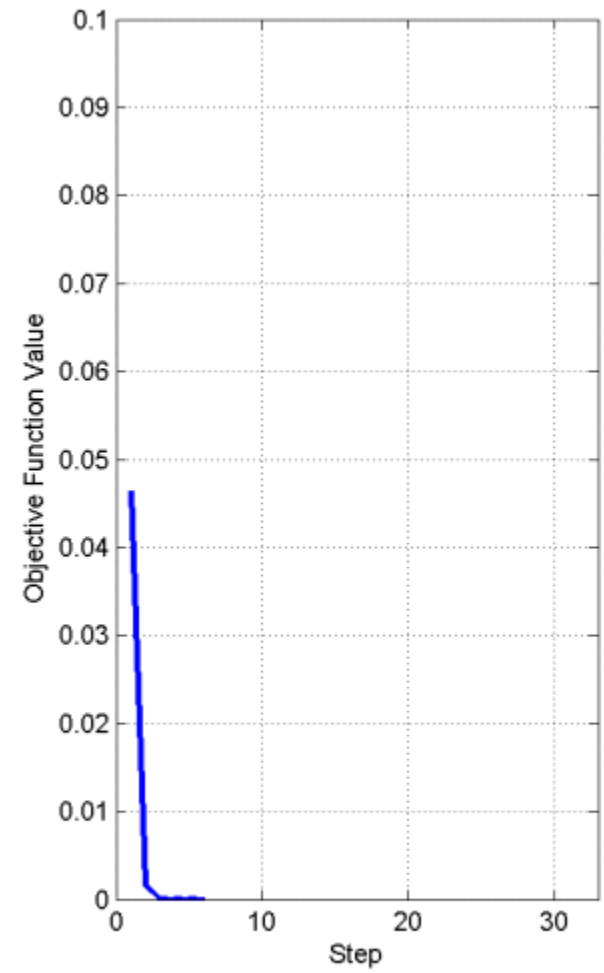
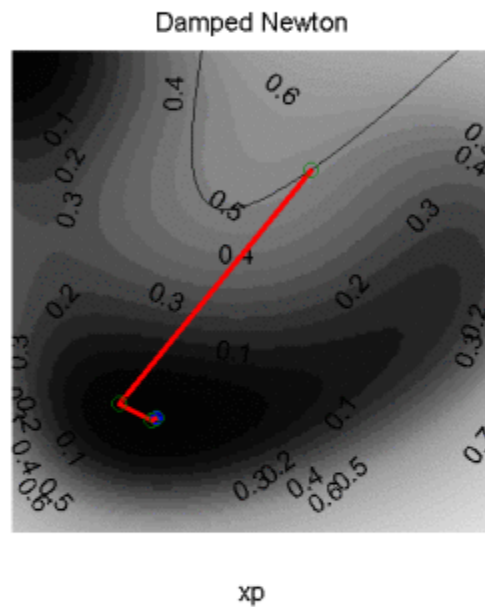
Gradient' vs. Newton's Method



yp



yp



Newton's Method for Logistic Regression

- Optimization Problem

$$\arg \min \frac{1}{N} \sum_{i=1}^N -y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))$$

- Gradient and Hessian Matrix

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) \mathbf{x}^{(i)}$$

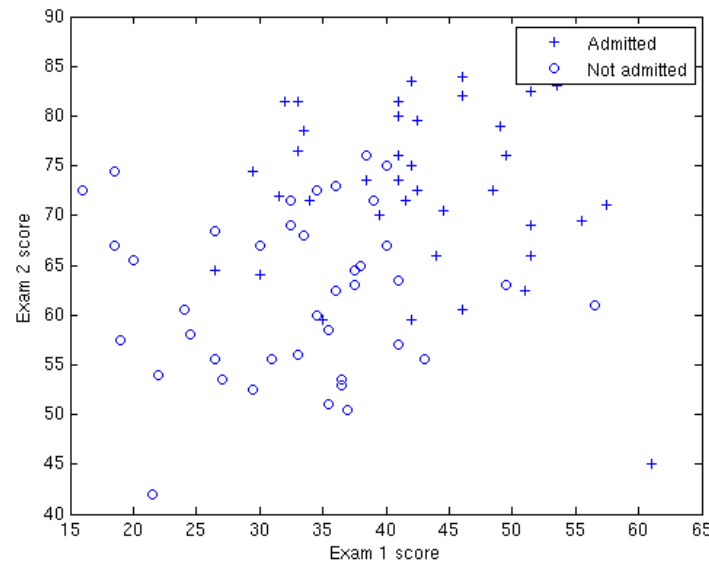
$$\mathbf{H} = \frac{1}{N} \sum_{i=1}^N h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})^T \left(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T$$

- Weight updating using Newton's method

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \mathbf{H}^{-1} \nabla J(\boldsymbol{\theta}^{(t)})$$

Practice 2: Logistic Regression

- Given the following training data:



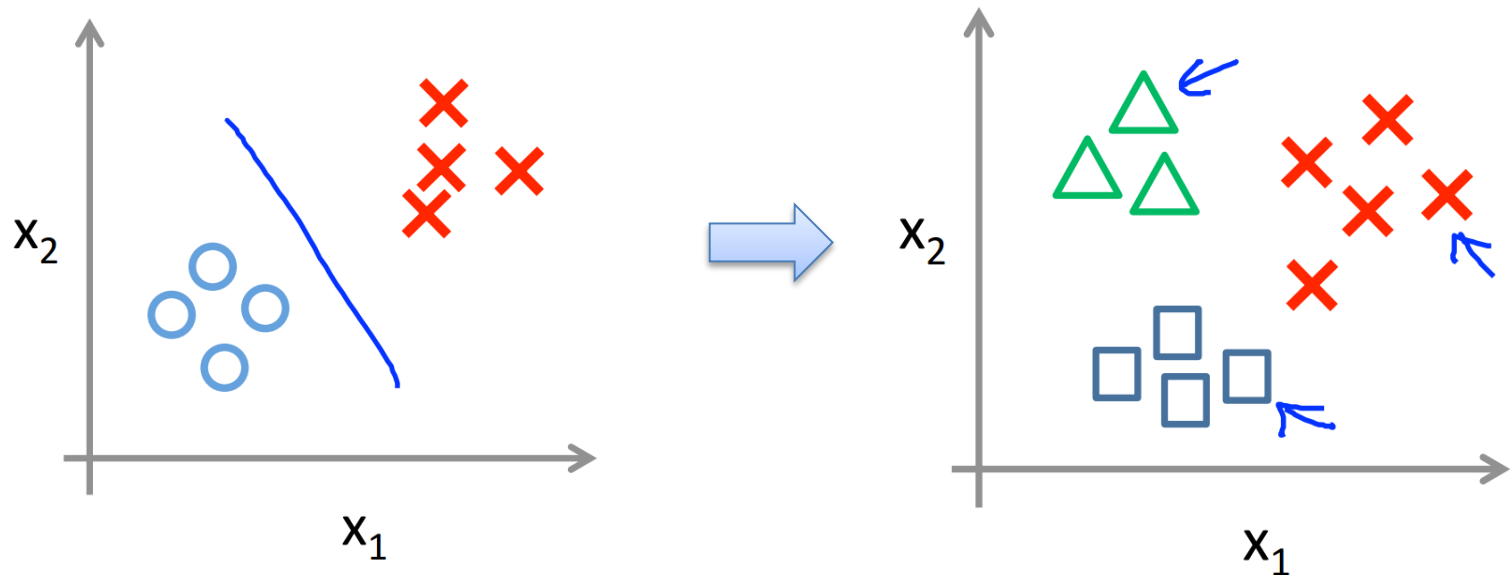
<http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=DeepLearning&doc=exercises/ex4/ex4.html>

- Implement 1) GD; 2) SGD; 3) Newton's Method for logistic regression, starting with the initial parameter $\theta = 0$.
- Determine how many iterations to use, and calculate for each iteration and plot your results.

Softmax Regression

Softmax Regression

- Softmax Regression is a multi-class classification model, also called Multi-class Logistic Regression;
- It is also known as the Maximum Entropy Model (in NLP);
- It is one of the most used classification algorithms.



Model Hypothesis

- Hypothesis

$$p(y = j | \mathbf{x}; \boldsymbol{\theta}) = h_j(\mathbf{x}) = \frac{e^{\boldsymbol{\theta}_j^T \mathbf{x}}}{1 + \sum_{j'=1}^{C-1} e^{\boldsymbol{\theta}_{j'}^T \mathbf{x}}}, j = 1, \dots, C-1$$

$$p(y = C | \mathbf{x}; \boldsymbol{\theta}) = h_C(\mathbf{x}) = \frac{1}{1 + \sum_{j'=1}^{C-1} \exp\{\boldsymbol{\theta}_{j'}^T \mathbf{x}\}}$$

- Hypothesis (Compact Form)

$$p(y = j | \mathbf{x}; \boldsymbol{\theta}) = h_j(\mathbf{x}) = \frac{e^{\boldsymbol{\theta}_j^T \mathbf{x}}}{\sum_{j'=1}^C e^{\boldsymbol{\theta}_{j'}^T \mathbf{x}}}, j = 1, 2, \dots, C, \text{ where } \boldsymbol{\theta}_C = \vec{0}$$

- Parameters

$$\boldsymbol{\theta}_{C \times M}$$

Maximum Likelihood Estimation

- (Conditional) Log-likelihood

$$\begin{aligned}l(\boldsymbol{\theta}) &= \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) && \text{Softmax Regression} \\&= \sum_{i=1}^N \log \prod_{j=1}^C \left(\frac{e^{\boldsymbol{\theta}_j^T \mathbf{x}}}{\sum_{j'=1}^C e^{\boldsymbol{\theta}_{j'}^T \mathbf{x}}} \right)^{1\{y^{(i)}=j\}} \\&= \sum_{i=1}^N \sum_{j=1}^C 1\{y^{(i)} = j\} \log \left(\frac{e^{\boldsymbol{\theta}_j^T \mathbf{x}}}{\sum_{j'=1}^C e^{\boldsymbol{\theta}_{j'}^T \mathbf{x}}} \right) \\&= \sum_{i=1}^N \sum_{j=1}^C 1\{y^{(i)} = j\} \log h_j(\mathbf{x}^{(i)})\end{aligned}$$

$$l(\boldsymbol{\theta}) = \sum_{i=1}^N y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \quad \text{Logistic Regression}$$

Gradient Descent Optimization

- Gradient

$$\frac{\partial \log h_j(\mathbf{x})}{\partial \boldsymbol{\theta}_k} = \begin{cases} (1 - h_k(\mathbf{x}))\mathbf{x}, & j = k \\ -h_k(\mathbf{x})\mathbf{x}, & j \neq k \end{cases}$$

$$\begin{aligned} \frac{\partial \sum_{j=1}^C 1\{y = j\} \log h_j(\mathbf{x})}{\partial \boldsymbol{\theta}_k} &= \begin{cases} (1 - h_k(\mathbf{x}))\mathbf{x}, & y = k \\ -h_k(\mathbf{x})\mathbf{x}, & y \neq k \end{cases} \\ &= (1\{y = k\} - h_k(\mathbf{x}))\mathbf{x} \end{aligned}$$

$$\frac{\partial l(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_k} = \sum_{i=1}^N \boxed{(1\{y^{(i)} = k\} - h_k(\mathbf{x}^{(i)})) \mathbf{x}^{(i)}}$$

Error × Feature

Gradient Descent Optimization

- Gradient Descent

$$\boldsymbol{\theta}_k := \boldsymbol{\theta}_k + \alpha \sum_{i=1}^N (1\{y^{(i)} = k\} - h_k(\mathbf{x}^{(i)})) \mathbf{x}^{(i)}$$

$$\text{where } h_k(\mathbf{x}) = \frac{e^{\boldsymbol{\theta}_k^T \mathbf{x}}}{\sum_{k'=1}^C e^{\boldsymbol{\theta}_{k'}^T \mathbf{x}}}, k = 1, 2, \dots, C$$

- Stochastic Gradient Descent

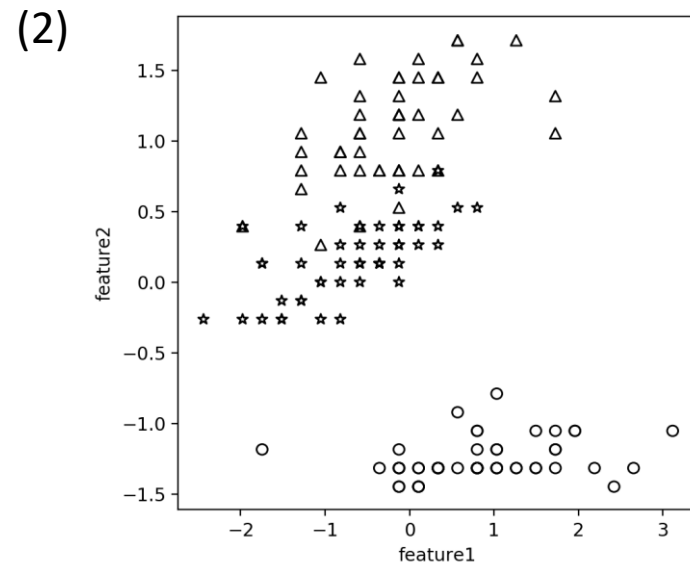
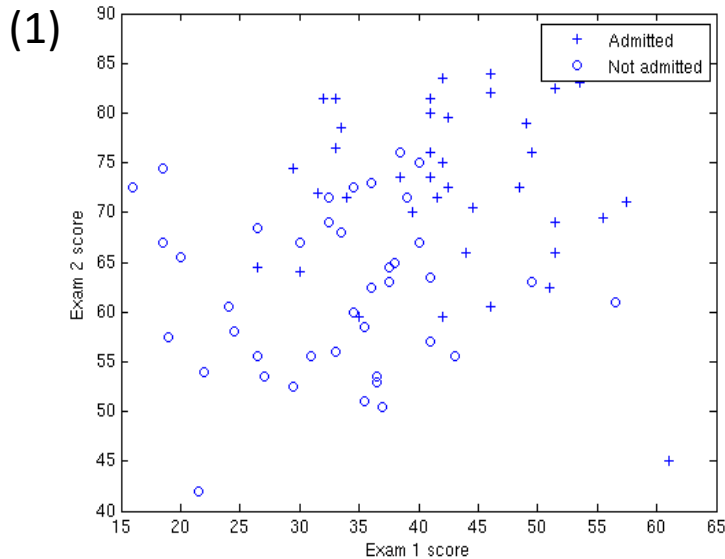
$$\boldsymbol{\theta}_k := \boldsymbol{\theta}_k + \alpha (1\{y = k\} - h_k(\mathbf{x})) \mathbf{x}$$

The other optimization methods

- Newton Method
- Quasi-Newton Method (BFGS)
- Limited Memory BFGS (L-BFGS)
- Conjugate Gradient
- GIS
- IIS
- ...

Practice 3: Softmax Regression

- Given the following training data sets:



- (1) <http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=DeepLearning&doc=exercises/ex4/ex4.html>
(2) <https://pan.baidu.com/s/1gU81bKslj8cRokOYEK1Jzw> password: w2a8

- For data set (1), implement logistic regression and softmax regression with 1) GD; 2) SGD.
- For data set (2), implement softmax regression with 1) GD; 2) SGD.
- Compare logistic regression and softmax regression.



Any Questions?