

Multiclass Logistic Regression

Sargur N. Srihari

University at Buffalo, State University of New York
USA

Topics in Linear Classification using Probabilistic Discriminative Models

- Generative vs Discriminative
 1. Fixed basis functions in linear classification
 2. Logistic Regression (two-class)
 3. Iterative Reweighted Least Squares (IRLS)
 4. **Multiclass Logistic Regression**
 5. Probit Regression
 6. Canonical Link Functions

Topics in Multiclass Logistic Regression

- Multiclass Classification Problem
- Softmax Regression
- Softmax Regression Implementation
- Softmax and Training
- One-hot vector representation
- Objective function and gradient
- Summary of concepts in Logistic Regression
- Example of 3-class Logistic Regression

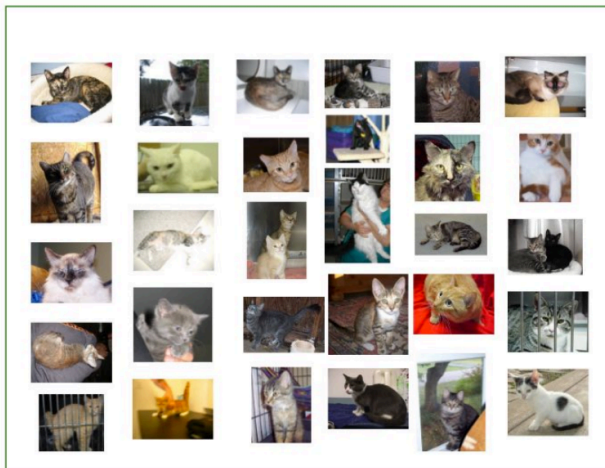
Multi-class Classification problem

Categories
 $K=10$

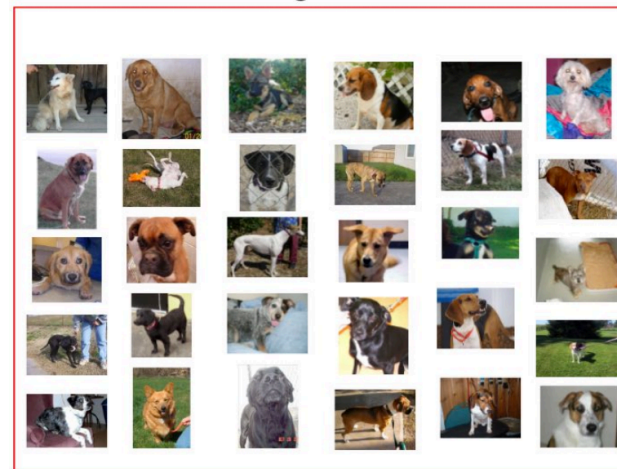
Examples
 $N=100$



Cats



Dogs



Softmax Regression

- In the two-class case $p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T\phi + b)$
 where $\phi = [\phi_1, \dots, \phi_M]^T$, $\mathbf{w} = [w_1, \dots, w_M]^T$ and $a = \mathbf{w}^T\phi + b$ is the “activation”
- For K classes, we work with soft-max function instead of logistic sigmoid (*Softmax regression*)

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where $a_k = \mathbf{w}_k^T \phi + b_k$, $k = 1, \dots, K$

$\mathbf{w}_k = [w_{k1}, \dots, w_{kM}]^T$ and $\mathbf{a} = \{a_1, \dots, a_K\}$

- We learn a set of K weight vectors $\{\mathbf{w}_1, \dots, \mathbf{w}_K\}$ and biases \mathbf{b}
- Arranging weight vectors as a matrix W

$$\mathbf{a} = W^T \phi + \mathbf{b}$$

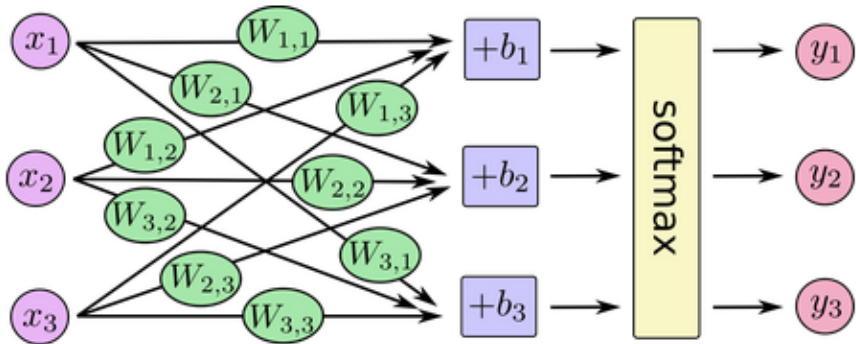
$$W = \begin{bmatrix} \mathbf{w}_1 \\ \cdot \\ \cdot \\ \mathbf{w}_K \end{bmatrix} = \begin{bmatrix} w_{11} & w_{1M} \\ \cdot & \cdot \\ \cdot & \cdot \\ w_{K1} & w_{KM} \end{bmatrix}$$

$$\mathbf{y} = \text{softmax}(\mathbf{a})$$

$$y_i = \frac{\exp(a_i)}{\sum_{j=1}^3 \exp(a_j)}$$

Softmax Regression Implementation

3-class Logistic Regression with 3 inputs



$$\mathbf{a} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\mathbf{W} = \begin{bmatrix} W_1 \\ W_2 \\ W_3 \end{bmatrix} = \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix}$$

$$\mathbf{y} = \text{softmax}(\mathbf{a})$$

$$y_i = \frac{\exp(a_i)}{\sum_{j=1}^3 \exp(a_j)}$$

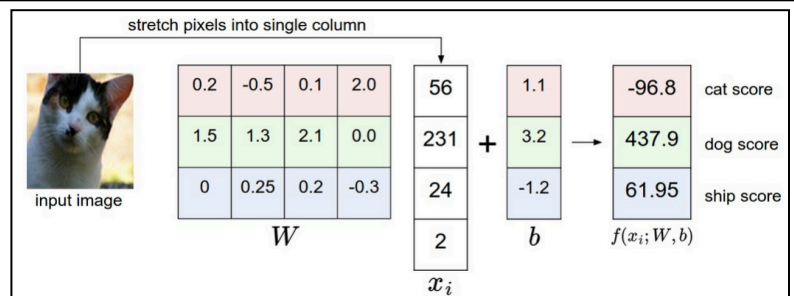
Network Computes

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \text{softmax} \begin{bmatrix} W_{1,1}x_1 + W_{1,2}x_2 + W_{1,3}x_3 + b_1 \\ W_{2,1}x_1 + W_{2,2}x_2 + W_{2,3}x_3 + b_2 \\ W_{3,1}x_1 + W_{3,2}x_2 + W_{3,3}x_3 + b_3 \end{bmatrix}$$

In matrix multiplication notation

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \text{softmax} \left(\begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \right)$$

An example



Softmax and Training

- We use maximum likelihood to determine the parameters $\{\mathbf{w}_k\}$, $k=1,\dots,K$

- The exp within softmax $\text{softmax}(\mathbf{a})_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$ works very well when training using log-likelihood

– Log-likelihood can undo the exp of softmax

$$\log \text{softmax}(\mathbf{a})_i = a_i - \log \sum_j \exp(a_j)$$

– Input a_i always has a direct contribution to cost

- Because this term cannot saturate, learning can proceed even if second term becomes very small

– First term encourages a_i to be pushed up

– Second term encourages all \mathbf{a} to be pushed down

Derivatives

- The multiclass logistic regression model is

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

- For maximum likelihood we will need the derivatives of y_k wrt all of the activations a_j
- These are given by

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$

– where I_{kj} are the elements of the identity matrix

One-hot vector representation

- Classes C_1, \dots, C_K represented by 1-of- K scheme

– One-hot vector:

- class C_k is a K -dim vector or $[t_1, \dots, t_K]^T$, $t_i \in \{0, 1\}$

– With $K=6$, class C_3 is $(0, 0, 1, 0, 0, 0)^T$ with $t_1=t_2=t_4=t_5=t_6=0$, & $t_3=1$

– The class probabilities obey

$$\sum_{k=1}^K p(C_k) = \sum_{k=1}^K t_k = 1$$

– If $p(t_k=1) = \mu_k$ then

$$p(C_k) = \prod_{k=1}^K \mu_k^{t_k} \text{ where } \boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$$

e.g., probability of C_3 is

$$p([0, 0, 1, 0, 0, 0]) = \mu_3$$

Why use one-hot representation?

If we used numerical categories 1,2,3...we would impute ordinality.

We can now use simpler Bernoulli instead of multinoulli

Target Matrix, T

- Classes have values $1, \dots, K$
- Each represented as a K -dimensional binary vector
- We have N labeled samples
 - So instead of target vector \mathbf{t} we have a target matrix T

$$T = \begin{array}{c} \xrightarrow{\text{Classes}} \\ \left[\begin{array}{cccc} t_{11} & \cdot & \cdot & t_{1K} \\ \cdot & & & \\ \cdot & & & \\ t_{N1} & & & t_{NK} \end{array} \right] \downarrow \text{Samples} \end{array}$$

Note that t_{nk} corresponds to sample n and class k

Objective Function & Gradient

- Likelihood of observations

$$T = \begin{bmatrix} t_{11} & \cdot & \cdot & t_{1K} \\ \cdot & & & \\ \cdot & & & \\ t_{N1} & & & t_{NK} \end{bmatrix} \text{ is}$$

$$p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

- Where, for feature vector ϕ_n $y_{nk} = y_k(\phi_n) = \frac{\exp(\mathbf{w}_k^T \phi_n)}{\sum_j \exp(\mathbf{w}_j^T \phi_n)}$

- Objective Function: negative log-likelihood

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

– Known as cross-entropy error for multi-class

- Gradient of error function wrt parameter \mathbf{w}_j

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \underbrace{(y_{nj} - t_{nj})}_{\text{Error}} \phi_n$$

$$\text{using } \sum_k t_{nk} = 1$$

Error x Feature Vector

$$y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \mathbf{w}_k^T \phi$$

$$\frac{\partial y_k}{\partial a_j} = y_k (I_{kj} - y_j) \text{ where } I_{kj} \text{ are elements of}$$

the identity matrix

Gradient Descent

- Has same form for gradient as for the sum of squares error function with the linear model and cross-entropy error for the logistic regression model
- i.e., product of the error $(y_{nj} - t_{nj})$ times the basis function ϕ_n
- We can use the sequential algorithm in which inputs are presented one at a time in which the weight vector is updated using

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_n$$

Newton-Raphson update gives IRLS

- Hessian matrix comprises blocks of size $M \times M$
 - Block j,k is given by

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T$$

- No of blocks is also $M \times M$, each corresponding to a pair of classes (with redundancy)
- Hessian matrix is positive-definite, therefore error function has a unique minimum
- Batch Algorithm based on Newton-Raphson

Summary of Logistic Regression concepts

- Definition of gradient and Hessian
- Gradient and Hessian in Linear Regression
- Gradient and Hessian in 2-class Logistic Regression

Definitions of Gradient and Hessian

- First derivative of a scalar function $E(\mathbf{w})$ with respect to a vector $\mathbf{w}=[w_1, w_2]^T$ is a vector called the *Gradient* of $E(\mathbf{w})$

$$\nabla E(\mathbf{w}) = \frac{d}{d\mathbf{w}} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_2} \end{bmatrix}$$

If there are M elements in the vector then Gradient is a $M \times 1$ vector

- Second derivative of $E(\mathbf{w})$ is a matrix called the *Hessian*

$$H = \nabla \nabla E(\mathbf{w}) = \frac{d^2}{d\mathbf{w}^2} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 E}{\partial w_1^2} & \frac{\partial^2 E}{\partial w_1 \partial w_2} \\ \frac{\partial^2 E}{\partial w_2 \partial w_1} & \frac{\partial^2 E}{\partial w_2^2} \end{bmatrix}$$

Hessian is a matrix with M^2 elements

- *Jacobian* matrix consists of first derivatives of a vector-valued function wrt a vector

Use of Gradient & Hessian in ML

Training samples: $n=1,..N$

Inputs: $M \times 1$ vectors

$$\phi(\mathbf{x}_n) = \begin{pmatrix} \phi_0(\mathbf{x}_n) & \phi_1(\mathbf{x}_n) & \dots & \phi_{M-1}(\mathbf{x}_n) \end{pmatrix}^T$$

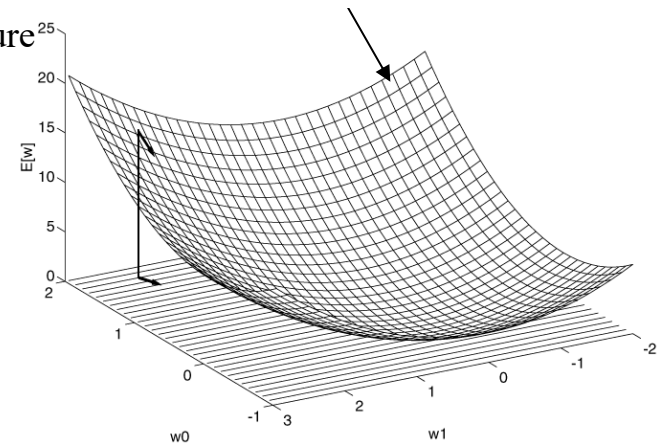
Outputs

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$\Phi_0(\mathbf{x})=1$, dummy feature

Error surface for $M=2$
a paraboloid with
a single global minimum



For Linear Regression (sum-of-squared error):

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [\mathbf{w}^T \phi(\mathbf{x}_n) - t_n]^2 \quad \text{where } \mathbf{w} = (w_0, w_1, \dots, w_{M-1})^T$$

For Logistic Regression (cross-entropy error):

$$E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\} \quad \text{where } y_n = \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))$$

For Stochastic Gradient Descent we need $\nabla E_n(\mathbf{w})$ where $E(\mathbf{w}) = \sum_n E_n(\mathbf{w})$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w})$$

For Newton-Raphson update we need both $\nabla E(\mathbf{w})$ and $H = \nabla \nabla E(\mathbf{w})$

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$$

Gradient and Hessian for Linear Regression

- Sum-of squared Errors (equivalent to maximum likelihood)

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [\mathbf{w}^T \phi(\mathbf{x}_n) - t_n]^2$$

$$\phi(\mathbf{x}_n) = \begin{pmatrix} \phi_0(\mathbf{x}_n) & \phi_1(\mathbf{x}_n) & \dots & \phi_{M-1}(\mathbf{x}_n) \end{pmatrix}^T$$

$$\mathbf{t} = (t_1, \dots, t_N)$$

- Gradient of E

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N [\mathbf{w}^T \phi_n - t_n] \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- Hessian of E

$$H = \nabla \nabla E(\mathbf{w}) = \Phi^T \Phi$$

- Newton-Raphson

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - (\Phi^T \Phi)^{-1} \{ \Phi^T \Phi \mathbf{w}^{(\text{old})} - \Phi^T \mathbf{t} \}$$

$$= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Which is the same solution as with Gradient Descent

Gradient & Hessian: 2-class Logistic Regression

- Cross-Entropy Error

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\} \text{ where } y_n = \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))$$

- Gradient of E

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi(\mathbf{x}_n) = \Phi^T (\mathbf{y} - \mathbf{t})$$

$$\mathbf{y} = (y_1, \dots, y_N)^T$$

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

- Hessian of E

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi(\mathbf{x}_n) \phi^T(\mathbf{x}_n) = \Phi^T \mathbf{R} \Phi$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

\mathbf{R} is $N \times N$ diagonal matrix with elements

$$R_{nn} = y_n(1 - y_n) = \mathbf{w}^T \phi(\mathbf{x}_n) (1 - \mathbf{w}^T \phi(\mathbf{x}_n))$$

$$\phi(\mathbf{x}_n) = \left(\phi_0(\mathbf{x}_n) \quad \phi_1(\mathbf{x}_n) \quad \dots \quad \phi_{M-1}(\mathbf{x}_n) \right)^T$$

Hessian is not constant and depends on \mathbf{w} through \mathbf{R}

Since \mathbf{H} is positive-definite (i.e., for arbitrary \mathbf{u} , $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$)

error function is a concave function of \mathbf{w} and so has a unique minimum

Gradient & Hessian: Multi-class Logistic Regression

- Cross-Entropy Error

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

$$T = \begin{bmatrix} t_{11} & \cdot & \cdot & t_{1K} \\ \cdot & & & \\ \cdot & & & \\ t_{N1} & & & t_{NK} \end{bmatrix}$$

- Gradient of E

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi(\mathbf{x}_n)$$

$$y_{nk} = y_k(\phi(\mathbf{x}_n)) = \frac{\exp(\mathbf{w}_k^T \phi(\mathbf{x}_n))}{\sum_j \exp(\mathbf{w}_j^T \phi(\mathbf{x}_n))}$$

- Hessian of E

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi(\mathbf{x}_n) \phi^T(\mathbf{x}_n)$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & & & \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & & & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$$\phi(\mathbf{x}_n) = \left(\phi_0(\mathbf{x}_n) \quad \phi_1(\mathbf{x}_n) \quad \dots \quad \phi_{M-1}(\mathbf{x}_n) \right)^T$$

Each element of the Hessian needs M multiplications and additions
 Since there are M^2 elements in the matrix the computation is $O(M^3)$

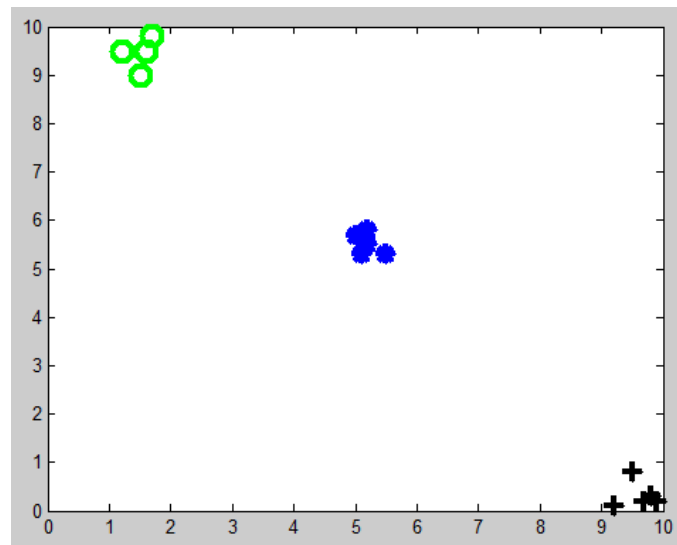
An Example of 3-class Logistic Regression

- Input Data

C1 =	
1.7000	9.8000
1.2000	9.5000
1.5000	9.0000
1.6000	9.5000
1.5000	9.0000

C2 =	
5.1000	5.3000
5.2000	5.5000
5.5000	5.3000
5.2000	5.8000
5.0000	5.7000

C3 =	
9.7000	0.2000
9.8000	0.3000
9.2000	0.1000
9.5000	0.8000
9.9000	0.2000



$\Phi_0(x)=1$, dummy feature

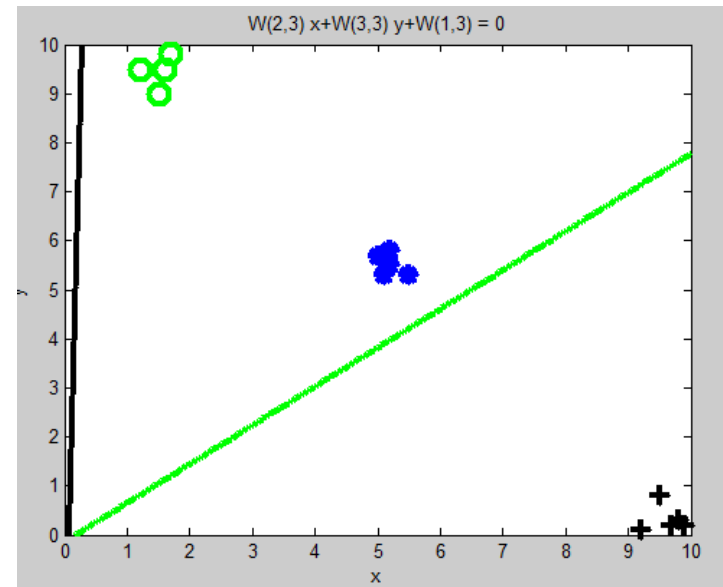
Three-class Logistic Regression

- Three weight vectors (Initial)

W =			
	0.4709	0.5486	0.1839
	-2.5932	6.0889	-2.3215
	3.2777	-1.2836	0.0483

- Gradient

Delta_E =			
	-0.0000	-0.0000	-0.0000
	19.7000	1.2000	-20.9000
	-21.4667	-2.2667	23.7333



- Hessian (9x9 with some 3 x 3 blocks repeated)

H =									
	3.3333	18.1333	16.8889	-1.6667	-9.0667	-8.4444	-1.6667	-9.0667	-8.4444
	18.1333	135.5022	50.9067	-9.0667	-67.7511	-25.4533	-9.0667	-67.7511	-25.4533
	16.8889	50.9067	131.5378	-8.4444	-25.4533	-65.7689	-8.4444	-25.4533	-65.7689
	-1.6667	-9.0667	-8.4444	3.3333	18.1333	16.8889	-1.6667	-9.0667	-8.4444
	-9.0667	-67.7511	-25.4533	18.1333	135.5022	50.9067	-9.0667	-67.7511	-25.4533
	-8.4444	-25.4533	-65.7689	16.8889	50.9067	131.5378	-8.4444	-25.4533	-65.7689
	-1.6667	-9.0667	-8.4444	-1.6667	-9.0667	-8.4444	3.3333	18.1333	16.8889
	-9.0667	-67.7511	-25.4533	-9.0667	-67.7511	-25.4533	18.1333	135.5022	50.9067
	-8.4444	-25.4533	-65.7689	-8.4444	-25.4533	-65.7689	16.8889	50.9067	131.5378

Final Weight Vector, Gradient and Hessian (3-class)

- Weight Vector

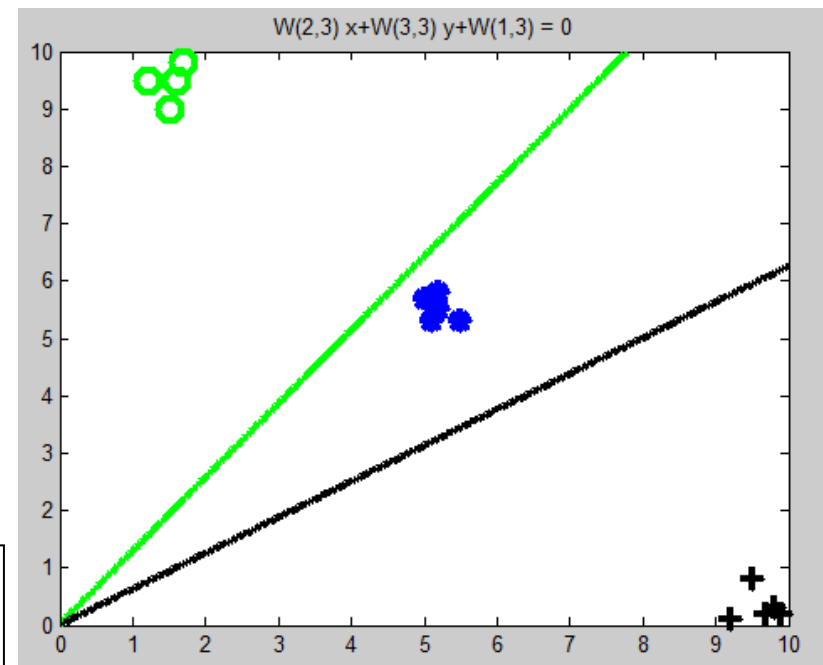
W =		
0.4685	0.7553	0.1523
-1.4499	1.0309	1.9186
1.8775	0.7173	-1.3135

- Gradient

Delta_E =		
0.0016	-0.0020	0.0004
0.0603	-0.0174	-0.0429
-0.0446	-0.0089	0.0535

- Hessian

H =								
0.0134	0.0459	0.0981	-0.0134	-0.0459	-0.0980	-0.0000	-0.0000	-0.0000
0.0459	0.2009	0.2924	-0.0459	-0.2007	-0.2923	-0.0000	-0.0001	-0.0001
0.0981	0.2924	0.7655	-0.0980	-0.2923	-0.7654	-0.0000	-0.0001	-0.0002
-0.0134	-0.0459	-0.0980	0.0224	0.1126	0.1251	-0.0090	-0.0668	-0.0270
-0.0459	-0.2007	-0.2923	0.1126	0.7359	0.4437	-0.0668	-0.5352	-0.1514
-0.0980	-0.2923	-0.7654	0.1251	0.4437	0.9018	-0.0270	-0.1514	-0.1365
-0.0000	-0.0000	-0.0000	-0.0090	-0.0668	-0.0270	0.0090	0.0668	0.0270
-0.0000	-0.0001	-0.0001	-0.0668	-0.5352	-0.1514	0.0668	0.5353	0.1516
-0.0000	-0.0001	-0.0002	-0.0270	-0.1514	-0.1365	0.0270	0.1516	0.1366



Number of iterations : 6

Error (Initial and Final): 38.9394, 1.5000e-009