

Iterative Reweighted Least Squares

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

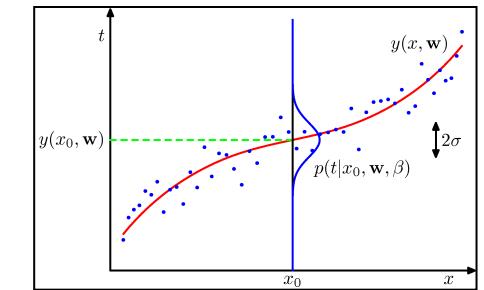
Topics in IRLS

- Linear and Logistic Regression
- Newton-Raphson Method
- Hessian
- IRLS for Linear Regression
- IRLS for Logistic Regression
- Numerical Example
- Hessian in Deep Learning

Recall Linear Regression

- Assuming Gaussian noise $p(t|x,w,\beta) = N(t | y(x,w), \beta^{-1})$
- And a linear model
- Likelihood has the form

$$y(x,w) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \phi(x)$$



$$p(t | X, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

- Log-likelihood becomes

$$\ln p(t | w, \beta) = \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \beta E_D(w)$$

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w^T \phi(x_n) \right\}^2$$

- Which is quadratic in input
- Derivative has the form
- Setting equal to zero and solving we get

$$\nabla \ln p(t | w, \beta) = \beta \sum_{n=1}^N \left\{ t_n - w^T \phi(x_n) \right\} \phi(x_n)^T$$

$$w_{ML} = \Phi^+ t$$

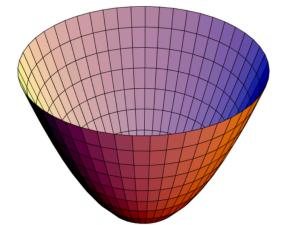
$$\Phi^+ = (\Phi^T \Phi)^{-1} \Phi^T$$

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & & & \\ \vdots & & & \\ \phi_0(x_N) & & & \phi_{M-1}(x_N) \end{pmatrix}$$

Linear and Logistic Regression

- In linear regression there is a closed-form max likelihood solution for determining w
 - on the assumption of Gaussian noise model
 - Due to quadratic dependence of log-likelihood on w

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w^T \phi(x_n) \right\}^2$$



- For logistic regression: No closed-form maximum likelihood solution
 - Due to nonlinearity of logistic sigmoid
- But departure from quadratic is not substantial
 - Error function is concave, i.e., unique minimum

$$E(w) = -\ln p(t | w) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

$$y_n = \sigma(w^T \phi_n)$$

What is IRLS?

- An iterative method to find solution \mathbf{w}^*
 - for linear regression and logistic regression
 - assuming least squares objective
- While simple gradient descent has the form

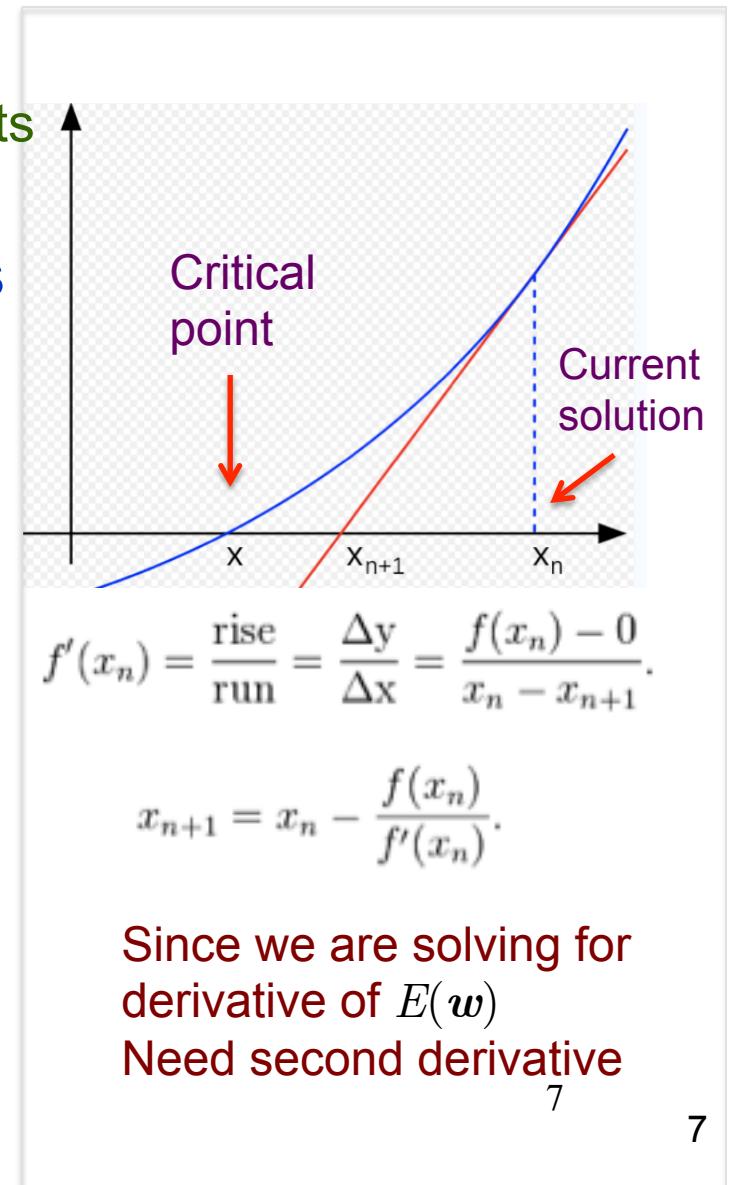
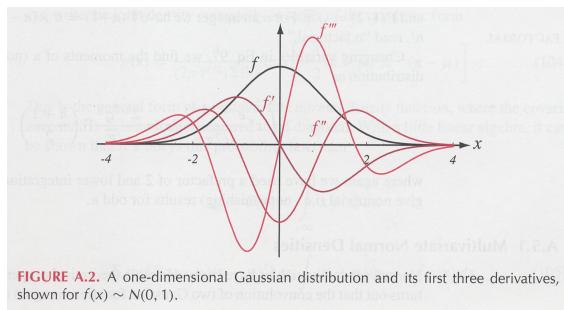
$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \eta \nabla E(\mathbf{w})$$

- IRLS uses the second derivative and has the form
- It is derived from Newton-Raphson method
 - where H is the Hessian matrix whose elements are the second derivatives of $E(\mathbf{w})$ wrt \mathbf{w}

Newton-Raphson Method (1-D)

- Based on second derivatives
 - Derivative of a function at x is the slope of its tangent at that point
- Illustration of 2nd and higher derivatives
 - Derivatives of Gaussian $p(x) \sim N(0, \sigma^2)$

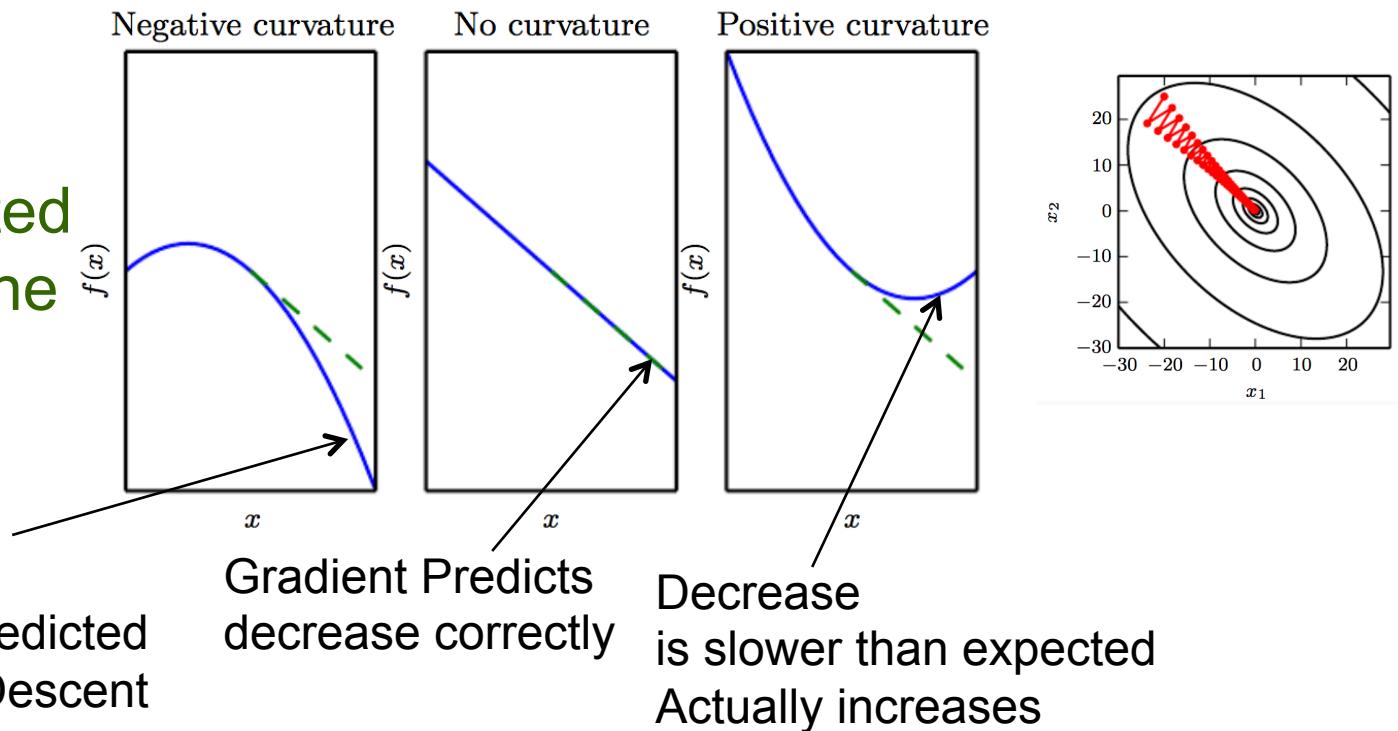
$$\begin{aligned}\frac{\partial}{\partial x} \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{-x}{\sqrt{2\pi}\sigma^3} e^{-x^2/(2\sigma^2)} = \frac{-x}{\sigma^2} p(x) \\ \frac{\partial^2}{\partial x^2} \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{1}{\sqrt{2\pi}\sigma^5} (-\sigma^2 + x^2) e^{-x^2/(2\sigma^2)} = \frac{-\sigma^2 + x^2}{\sigma^4} p(x) \\ \frac{\partial^3}{\partial x^3} \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} \right] &= \frac{1}{\sqrt{2\pi}\sigma^7} (3x\sigma^2 - x^3) e^{-x^2/(2\sigma^2)} = \frac{-3x\sigma^2 - x^3}{\sigma^6} p(x),\end{aligned}$$



Second derivative measures curvature

- Derivative of a derivative
- Quadratic functions with different curvatures

Dashed line is value of cost function predicted by gradient alone



Beyond Gradient: Jacobian and Hessian matrices

- Sometimes we need to find all derivatives of a function whose input and output are both vectors
- If we have function $f: R^m \rightarrow R^n$
 - Then the matrix of partial derivatives is known as the Jacobian matrix J defined as

$$J_{i,j} = \frac{\partial}{\partial x_j} f(x)_i$$

Second derivative

- Derivative of a derivative
- For a function $f: R^n \rightarrow R$ the derivative wrt x_i of the derivative of f wrt x_j is denoted as $\frac{\partial^2}{\partial x_i \partial x_j} f$
- In a single dimension we can denote $\frac{\partial^2}{\partial x^2} f$ by $f''(x)$
- Tells us how the first derivative will change as we vary the input
- This important as it tells us whether a gradient step will cause as much of an improvement as based on gradient alone

Hessian

- Second derivative with many dimensions
- $H(f)(\mathbf{x})$ is defined as
$$H(f)(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$$
- Hessian is the Jacobian of the gradient
- Hessian matrix is symmetric, i.e., $H_{i,j} = H_{j,i}$
 - anywhere that the second partial derivatives are continuous
 - So the Hessian matrix can be decomposed into a set of real eigenvalues and an orthogonal basis of eigenvectors
 - Eigenvalues of H are useful to determine learning rate as seen in next two slides

Role of eigenvalues of Hessian

- Second derivative in direction d is $d^T H d$
 - If d is an eigenvector, second derivative in that direction is given by its eigenvalue
 - For other directions, weighted average of eigenvalues (weights of 0 to 1, with eigenvectors with smallest angle with d receiving more value)
- Maximum eigenvalue determines maximum second derivative and minimum eigenvalue determines minimum second derivative

Learning rate from Hessian

- Taylor's series of $f(\mathbf{x})$ around current point $\mathbf{x}^{(0)}$

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^T \mathbf{g} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(0)})^T H(\mathbf{x} - \mathbf{x}^{(0)})$$

- where \mathbf{g} is the gradient and H is the Hessian at $\mathbf{x}^{(0)}$
 - If we use learning rate ε the new point \mathbf{x} is given by $\mathbf{x}^{(0)} - \varepsilon \mathbf{g}$. Thus we get

$$f(\mathbf{x}^{(0)} - \varepsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \varepsilon \mathbf{g}^T \mathbf{g} + \frac{1}{2} \varepsilon^2 \mathbf{g}^T H \mathbf{g}$$

- There are three terms:
 - original value of f ,
 - expected improvement due to slope, and
 - correction to be applied due to curvature
- Solving for step size when correction is least gives

$$\varepsilon^* \approx \frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T H \mathbf{g}}$$

Another 2nd Derivative Method

- Using Taylor's series of $f(\mathbf{x})$ around current $\mathbf{x}^{(0)}$

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^T \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^T H(f)(\mathbf{x} - \mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)})$$

- solve for the critical point of this function to give

$$\mathbf{x}^* = \mathbf{x}^{(0)} - H(f)(\mathbf{x}^{(0)})^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$$

- When f is a quadratic (positive definite) function use solution to jump to the minimum function directly
 - When not quadratic apply solution iteratively
- Can reach critical point much faster than gradient descent
 - But useful only when nearby point is a minimum

Two applications of IRLS

- IRLS is applicable to both Linear Regression and Logistic Regression
- We discuss both, for each we need

1. Model

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

$$p(C_1 | \Phi) = y(\Phi) = \sigma(\mathbf{w}^T \Phi)$$

2. Objective function $E(\mathbf{w})$

- Linear Regression: Sum of Squared Errors

- Logistic Regression: Bernoulli Likelihood Function

3. Gradient

$$\nabla E(\mathbf{w})$$

4. Hessian

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

5. Newton-Raphson update

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

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

$$p(t | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \left\{ 1 - y_n \right\}^{1-t_n}$$

IRLS for Linear Regression

1. Model:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

2. Error Function: Sum of Squares:

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

for data set $\mathbf{X} = \{\mathbf{x}_n, t_n\} \quad n=1,..N$

3. Gradient of Error Function is:

$$\begin{aligned} \nabla E(\mathbf{w}) &= \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}_n - t_n) \boldsymbol{\phi}_n \\ &= \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{t} \end{aligned}$$

where $\boldsymbol{\Phi}$ is the $N \times M$ design matrix whose n^{th} row is given by $\boldsymbol{\phi}_n^T$

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

4. Hessian is:

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

5. Newton-Raphson Update:

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

Substituting: $H = \boldsymbol{\Phi}^T \boldsymbol{\Phi}$ and $\nabla E(\mathbf{w}) = \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{t}$

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \{ \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w}^{(\text{old})} - \boldsymbol{\Phi}^T \mathbf{t} \} \\ &= (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t} \end{aligned}$$

which is the standard least squares solution

Since it is independent of \mathbf{w} , Newton-Raphson gives exact solution in one step

IRLS for Logistic Regression

1. Model: $p(C_1|\phi) = y(\phi) = \sigma(w^T\phi)$ Likelihood:

for data set $\{\phi_n, t_n\}$, $t_n \in \{0,1\}$, $y_n = \phi(x_n)$

$$p(t | w) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

2. Objective Function: Negative log-likelihood:

$$E(w) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

Cross-entropy error function

3. Gradient of Error Function is:

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

where Φ is the $N \times M$ design matrix whose n^{th} row is given by ϕ_n^T

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & & & \\ \vdots & & & \\ \phi_0(x_N) & & & \phi_{M-1}(x_N) \end{pmatrix}$$

4. Hessian is:

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

R is $N \times N$ diagonal matrix with elements

$$R_{nn} = y_n(1 - y_n) = w^T \phi_n (1 - w^T \phi_n)$$

Hessian is not constant and depends on w through R . Since H is positive-definite (i.e., for arbitrary u , $u^T H u > 0$) error function is a concave function of w and so has a unique minimum

5. Newton-Raphson Update:

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

Substituting $H = \Phi^T R \Phi$ and $\nabla E(w) = \Phi^T(y - t)$

$$\begin{aligned} w^{(\text{new})} &= w^{(\text{old})} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t) \\ &= (\Phi^T R \Phi)^{-1} \{ \Phi \Phi w^{(\text{old})} - \Phi^T (y - t) \} \\ &= (\Phi^T R \Phi)^{-1} \Phi^T R z \end{aligned}$$

where z is a N -dimensional vector with elements $z = \Phi w^{(\text{old})} - R^{-1} (y - t)$

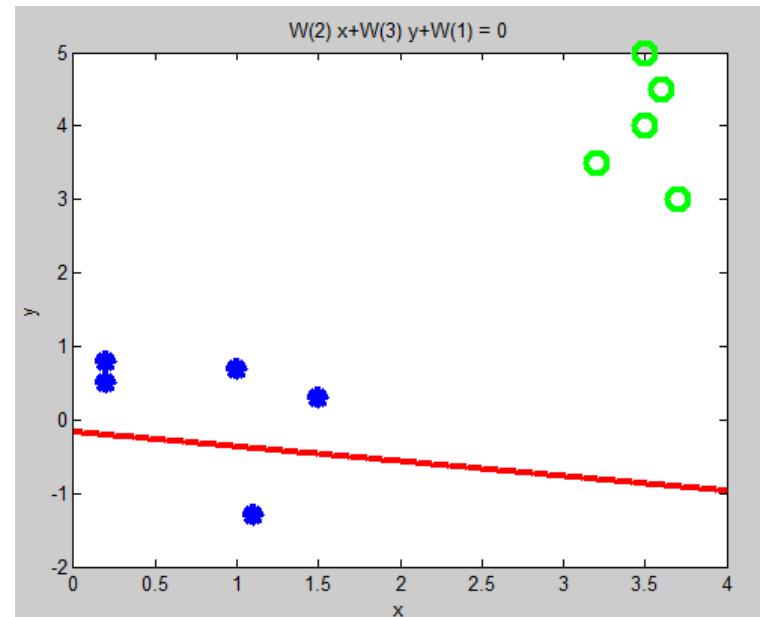
Update formula is a set of normal equations.

Since Hessian depends on w apply them iteratively each time using the new weight vector

Initial Weight Vector, Gradient and Hessian (2-class)

- Weight vector

```
w =
0.1117
0.1363
0.6787
```



- Gradient

```
Delta_E =
0
6.7500
9.5000
```

- Hessian

```
H =
3.5000    5.3750    5.2500
5.3750   17.4825   17.4950
5.2500   17.4950   22.4150
```

Final Weight Vector, Gradient and Hessian (2-class)

- Weight Vector

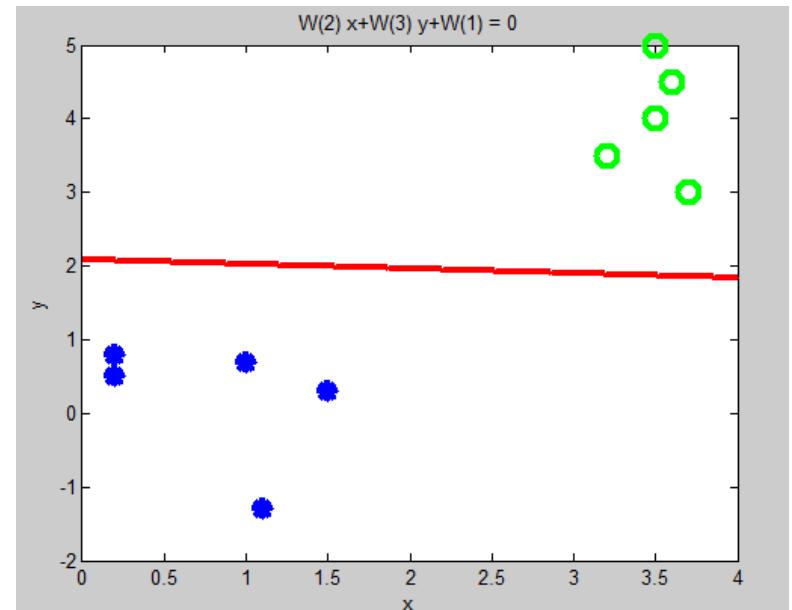
```
W =
704.5915
-20.9086
-337.6170
```

- Gradient

```
Delta_E =
-12.3917
-1.6321
4.9025
```

- Hessian

```
H =
1.0000 0.0000 0.0000
0.0000 1.0000 0.0000
0.0000 0.0000 1.0000
```



Number of iterations : 10

Error (Initial and Final): 15.0642, 1.0000e-009

Hessian in Deep Learning

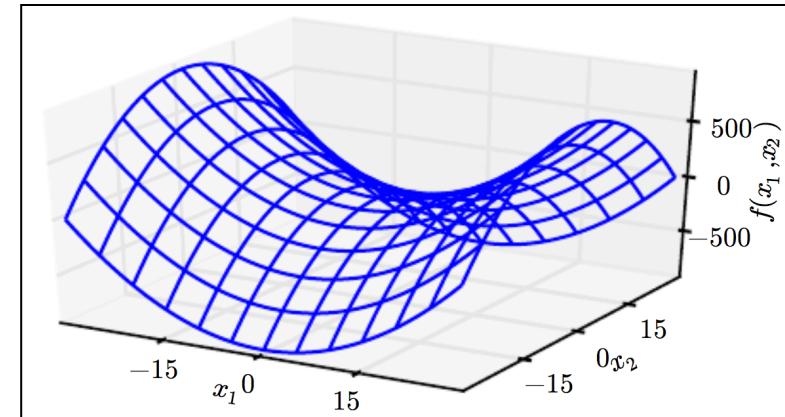
- The second derivative can be used to determine whether a critical point is a local minimum, maximum or a saddle point
 - On a critical point $f'(x)=0$
 - When $f''(x)>0$ the first derivative $f'(x)$ increases as we move to the right and decreases as we move left
 - We conclude that x is a local minimum
 - For local maximum, $f'(x)=0$ and $f''(x)<0$
 - When $f''(x)=0$ test is inconclusive: x may be a saddle point or part of a flat region

Multidimensional Second derivative test

- In multiple dimensions, we need to examine second derivatives of all dimensions
- Eigendecomposition generalizes the test
- Test eigenvalues of Hessian to determine whether critical point is a local maximum, local minimum or saddle point
- When H is positive definite (all eigenvalues are positive) the point is a local minimum
- Similarly negative definite implies a maximum

Saddle point

- Contains both positive and negative curvature
- Function is $f(x) = x_1^2 - x_2^2$



- Along axis x_1 , function curves upwards: this axis is an eigenvector of H and has a positive value
- Along x_2 , function curves downwards; its direction is an eigenvector of H with negative eigenvalue
- At a saddle point eigen values are both positive and negative